

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第5047988号

(P5047988)

(45) 発行日 平成24年10月10日(2012.10.10)

(24) 登録日 平成24年7月27日(2012.7.27)

(51) Int.Cl.

F I

G O 6 F 12/00 (2006.01)

G O 6 F 12/00 5 3 1 D

G O 6 F 3/06 (2006.01)

G O 6 F 3/06 3 0 1 F

G O 6 F 12/00 5 4 5 A

請求項の数 20 (全 89 頁)

(21) 出願番号 特願2008-548792 (P2008-548792)  
 (86) (22) 出願日 平成18年11月30日(2006.11.30)  
 (65) 公表番号 特表2009-522659 (P2009-522659A)  
 (43) 公表日 平成21年6月11日(2009.6.11)  
 (86) 国際出願番号 PCT/US2006/061431  
 (87) 国際公開番号 W02008/069811  
 (87) 国際公開日 平成20年6月12日(2008.6.12)  
 審査請求日 平成21年11月30日(2009.11.30)  
 (31) 優先権主張番号 60/754,726  
 (32) 優先日 平成17年12月29日(2005.12.29)  
 (33) 優先権主張国 米国(US)  
 (31) 優先権主張番号 11/371,304  
 (32) 優先日 平成18年3月8日(2006.3.8)  
 (33) 優先権主張国 米国(US)

(73) 特許権者 507303550  
 アマゾン・テクノロジー・インコーポ  
 レーテッド  
 アメリカ合衆国・89507・ネバダ州・  
 レノ・ピーオーボックス 8102  
 (74) 代理人 100064621  
 弁理士 山川 政樹  
 (74) 代理人 100098394  
 弁理士 山川 茂樹  
 (72) 発明者 バーミューレン, アラン・エイチ  
 アメリカ合衆国・98144・ワシントン  
 州・シアトル・12ティエイチ アベニュー  
 サウス・1200・スイート 1200

最終頁に続く

(54) 【発明の名称】 ウェブサービスクライアントインターフェースを有する分散型ストレージシステム

(57) 【特許請求の範囲】

【請求項1】

システムメモリに接続した1又は複数のプロセッサを備えたデータストレージ・ウェブ  
 サービス用のコンピュータにより実行可能な方法であって、当該方法は、

クライアントコンピュータに対してデータストレージ・ウェブサービスを提供するステ  
 ップであって、

前記データストレージ・ウェブサービス用のコンピュータは、

ウェブサービス要求を通じて前記クライアントコンピュータに対して適用可能である  
 データストレージのウェブサービス操作を定義するウェブサービスアプリケーションプロ  
 グラミングインターフェース(API)を実装し、

1つ以上の前記ウェブサービス操作を特定する前記ウェブサービスAPIに従ってフ  
 ォーマットされたウェブサービス要求を受信するために、インターネット系のアプリケー  
 ション層データ転送プロトコルに従ってアドレス可能であり、

データオブジェクトを格納するための前記ウェブサービスAPIに従ってフォーマッ  
 トされたウェブサービス要求を受信するのに応答して、前記クライアントコンピュータか  
 ら供給されたデータオブジェクトを格納する、ステップと、

前記インターネット系のアプリケーション層データ転送プロトコルに従って、データオ  
 ブジェクトにアクセスするためのクライアント要求を示す前記ウェブサービスAPIに従  
 ってフォーマットされたウェブサービス要求を受信するステップであって、前記データオ  
 ブジェクトの所定の1つにアクセスするための前記クライアント要求の所定の1つは、前

10

20

記所定データオブジェクトに対応するキー値を含む、ステップと、

複数のストレージノード上に前記データオブジェクトのコピーを保存するステップであって、前記コピーのそれぞれは、それぞれのロケータ値を介してアクセス可能であり、前記ロケータ値のそれぞれは、前記データストレージ・ウェブサービス内で固有である、ステップと、

前記データオブジェクトのそれぞれに対するそれぞれのキーマップエントリを保存するステップであって、前記任意データオブジェクトについては、前記それぞれのキーマップエントリは、前記クライアント特定キー値と、前記所定のデータオブジェクトのそれぞれの保存されたコピーに対応するそれぞれのロケータ値を含む、ステップと、

前記任意クライアント要求の受信に応じて、前記キー値に対応する1つまたは複数のロケータ値を識別するために、前記それぞれのキーマップエントリにアクセスするステップであって、前記1つまたは複数のロケータ値の特定の1つについては、対応するストレージノードにアクセスし、対応するコピーを取り出すステップと、  
を備えることを特徴とする方法。

【請求項2】

前記対応するコピーを取り出す前に、さらに前記所定のクライアント要求が、前記所定のデータオブジェクトへアクセスするのに十分な特権を有するかどうかを判断するステップと、前記所定のクライアント要求が十分な権利が与えられていない場合には、前記所定のクライアント要求を拒否するステップと、をさらに備えることを特徴とする請求項1に記載の方法。

【請求項3】

前記インターネット系のアプリケーション層データ転送プロトコルに従って、データオブジェクトを保存するために前記ウェブサービスAPIに従ってフォーマットされたクライアント要求を受信するステップをさらに備え、前記データオブジェクトの特定の1つを保存するための前記クライアント要求の特定の1つは、前記特定データに対応するキー値を含むことを特徴とする請求項1または2に記載の方法。

【請求項4】

前記特定データオブジェクトを保存するための料金を判断するステップをさらに備えたことを特徴とする請求項3に記載の方法。

【請求項5】

前記特定のクライアント要求に応じて、前記特定のデータオブジェクトの1つ以上のコピーを1つまたは複数の対応するストレージノードに保存するステップと、

前記特定のデータオブジェクトの所定のコピーの保存に応じて、前記所定のコピーに対応するロケータ値を受信ステップと、をさらに備えることを特徴とする請求項3に記載の方法。

【請求項6】

前記特定のデータオブジェクトの前記1つまたは複数のコピーが保存される、前記1つまたは複数の対応するストレージノードを、ストレージ規定に従って選択するステップをさらに備えたことを特徴とする請求項4に記載の方法。

【請求項7】

前記ストレージ規定は、前記特定のデータオブジェクトを保存するための前記特定のクライアント要求が完了することを示す前に、対応するストレージノードへ永久に保存されたことを示すために必要となる多くのコピーを特定することを特徴とする請求項6に記載の方法。

【請求項8】

前記ストレージ規定は、生成される前記特定のデータオブジェクトの所望のコピーの数をさらに特定することを特徴とする請求項7に記載の方法。

【請求項9】

前記キーマップエントリの所定の1つに対して、前記所定のキーマップエントリのそれぞれのロケータ値に対応する各コピーが、アクセス可能かどうかの判断をするために、前

10

20

30

40

50

記それぞれのキーマップエントリを調査するステップをさらに備えることを特徴とする請求項 8 に記載の方法。

【請求項 1 0】

前記所定のキーマップエントリについて、前記所定のキーマップエントリのそれぞれのロケータ値に対応するアクセス可能なコピーの数が、前記コピーの所望の数よりも少ない場合、前記コピーの所望の数を満たすために十分である付加的なコピーを作成するステップをさらに備えることを特徴とする請求項 9 に記載の方法。

【請求項 1 1】

前記複数のストレージノードは複数の領域に分散され、前記ストレージ規定は、前記特定のデータオブジェクトを保存するための前記特定のクライアント要求が完了することを示す前に、前記 1 つまたは複数のコピーが永久に保存されることを示すことが必要となる最少数の領域を特定することを特徴とする請求項 6 に記載の方法。

10

【請求項 1 2】

前記複数の領域の任意の 2 つの領域間において、ストレージノードの不具合の可能性の相関関係がしきい値未満であることを特徴とする請求項 1 1 に記載の方法。

【請求項 1 3】

前記ストレージ規定は、可能な場合、前記 1 つまたは複数のコピーのうちの少なくとも 1 つが、前記領域の所定のものに位置するストレージノードに書き込まれることをさらに特定することを特徴とする請求項 1 1 に記載の方法。

【請求項 1 4】

20

前記データストレージ・ウェブサービス用のコンピュータによって受信された前記所定のデータオブジェクトへのアクセスのための前記クライアント要求のうちの特定の 1 つは、前記所定のデータオブジェクトの特定のコピーに対応する特定のロケータ値を含むことを特徴とする請求項 1 1 に記載の方法。

【請求項 1 5】

アクセスのための前記特定のクライアント要求の受信に応じて、前記キーマップエントリにアクセスすることなく、前記特定のロケータ値を介し、対応するストレージノードから前記特定のコピーを取り出すステップをさらに備えることを特徴とする請求項 1 4 に記載の方法。

【請求項 1 6】

30

前記インターネット系のアプリケーション層データ転送プロトコルは、Representational State Transfer (REST) ウェブサービスモデルを実行し、前記インターネット系のアプリケーション層データ転送プロトコルに従って前記所定のクライアント要求を受信することは、Hyper text Transfer Protocol (HTTP) のバージョンに従いフォーマット化された要求を受信することを含むことを特徴とする請求項 1 に記載の方法。

【請求項 1 7】

前記インターネット系のアプリケーション層データ転送プロトコルは、文書ベースウェブサービスモデルを実行し、

前記インターネット系のアプリケーション層データ転送プロトコルに従って前記所定のクライアント要求を受信することは、Simple Object Access Protocol (SOAP) のバージョンに従い、カプセル化した文書の受信をすることを含み、前記所定のクライアント要求の内容は、前記文書に含まれ、XML のバージョンに従いフォーマットされることを特徴とする請求項 1 に記載の方法。

40

【請求項 1 8】

前記データオブジェクトのそれぞれに対するそれぞれのキーマップエントリを保存するステップは、階層的に配置された複数のインデックスノードを含み、それぞれが関連するタグ値を有するインデックスデータ構造内に保存されたキーマップエントリを索引付けするステップを含み、

前記キーマップエントリのそれぞれは、インデックスノードのそれぞれの 1 つと対応し

50

、所定の対応するインデックスノードを有する前記所定のキーマップエントリについては、前記所定の対応するインデックスノードの各上位と関連する各タグ値が、前記所定のキー値のプレフィックスであることを特徴とする請求項 1 に記載の方法。

【請求項 19】

前記請求項 1 乃至 18 のいずれかに記載の方法を実施するように構成された複数のコンピュータより構成されるシステム。

【請求項 20】

インストラクションを備えるコンピュータにアクセス可能な媒体であって、前記インストラクションは、データストレージ・ウェブサービス用のコンピュータにより実行され、クライアントコンピュータに対してデータストレージ・ウェブサービスを提供するステップであって、

前記データストレージ・ウェブサービス用のコンピュータは、

ウェブサービス要求を通じて前記クライアントコンピュータに対して適用可能であるデータストレージのウェブサービス操作を定義するウェブサービスアプリケーションプログラミングインターフェース (API) を実装し、

1 つ以上の前記ウェブサービス操作を特定する前記ウェブサービス API に従ってフォーマットされたウェブサービス要求を受信するために、インターネット系のアプリケーション層データ転送プロトコルに従ってアドレス可能であり、

データオブジェクトを格納するための前記ウェブサービス API に従ってフォーマットされたウェブサービス要求を受信するのに応答して、前記クライアントコンピュータから供給されたデータオブジェクトを格納する、ステップと、

データオブジェクトへアクセスするためのクライアント要求を示す前記ウェブサービス API に従ってフォーマットされたウェブサービス要求を処理するステップであって、前記データオブジェクトにアクセスするための前記クライアント要求は、前記インターネット系のアプリケーション層データ転送プロトコルに従って受信され、前記データオブジェクトの所定の 1 つへアクセスするための前記クライアント要求の所定の 1 つは、前記所定のデータオブジェクトに対応するキー値を含む、ステップと、

複数のストレージノード上に前記データオブジェクトのコピーを保存することを指示するステップであって前記ロケータ値のそれぞれは、前記データストレージ・ウェブサービス内で固有であるステップと、

前記データオブジェクトのそれぞれに対するそれぞれのキーマップエントリを保存することを指示するステップであって、前記所定のデータオブジェクトについては、前記それぞれのキーマップエントリは、前記クライアント特定キー値と、前記所定のデータオブジェクトのそれぞれの保存されたコピーに対応するそれぞれのロケータ値とを含むステップと、

を実行することが可能であり、

前記所定のクライアント要求を処理するステップは、前記キー値に対応する 1 つまたは複数のロケータ値を識別するための前記それぞれのキーマップエントリへアクセスするステップを含み、前記 1 つまたは複数の特定の 1 つについては、対応するストレージノードにアクセスし、対応するコピーを取り出すステップを含むことを特徴とするインストラクションを備えることを特徴とするコンピュータにアクセス可能な媒体。

【発明の詳細な説明】

【技術分野】

【0001】

本発明はデータストレージシステム、より具体的には、ウェブサービスとしてストレージへのアクセスを提供するように構成されるストレージシステムに関連する。

【背景技術】

【0002】

多くの異なるコンピュータアプリケーションは、様々な種類のアプリケーションデータの持続的保存のために、あるストレージ媒体に依存する。例えば、一般的なオフィスアプ

10

20

30

40

50

リケーションおよびマルチメディアアプリケーションは、中でも、文書、集計表、静止画像、音声及びビデオデータなどの様々な種類やフォーマットのアプリケーションデータを作成し、使用する。このようなデータは頻繁に、ユーザが繰り返しアクセス又は使用できるように保存される。例えば、ユーザは、一定期間にわたり多くの文書又は他のデータを保存し、作業することを希望する場合があります、必要な時に予想できる状態で、データを速やかに使用できることを望む場合がある。

#### 【 0 0 0 3 】

従来のコンピュータ関係システムにおいて、持続的なアプリケーションデータの保存のためにアプリケーションによって使用されるストレージ媒体は、光学又は半導体ストレージデバイスも使用されるが、最も一般的には磁気固定ドライブ又は「ハードドライブ」である。このようなデバイスは、アプリケーションを実行するコンピュータシステム内に統合されているか、あるいはローカルの周辺インターフェース又はネットワークを介してそのシステムにアクセス可能である。一般的に、アプリケーションストレージとしての役割をするデバイスは、ファイルシステムインターフェースなど、ストレージアクセスを必要とする様々なアプリケーションへの一貫したストレージインターフェースを提供するためのデバイスレベルでの動作を管理する、オペレーティングシステムによって管理される。

#### 【 0 0 0 4 】

この従来型モデルのアプリケーションストレージには、いくつかの限界が存在する。第1に、概して従来型モデルは、アプリケーションデータのアクセス性を制限する。例えば、アプリケーションデータが、特定のコンピュータシステムのローカルハードドライブに保存される場合、他のシステムで実行されているアプリケーションからアクセス不可能である。データがネットワークでアクセス可能なデバイスに保存される場合でも、周辺ネットワーク外のシステムで実行するアプリケーションは、そのデバイスにアクセスすることが不可能である場合がある。例えば、セキュリティ上の理由で、企業の外部のシステムが企業内のシステム又はリソース情報にアクセスできないように、企業は一般にローカルエリアネットワーク（LAN）へのアクセスを制限する。したがって、携帯デバイス上（例えば、ノート型又は手持ちサイズコンピュータ、個人情報端末、携帯電話デバイスなど）で実行するアプリケーションは、固定されたシステム又はネットワークに持続的に関連するデータへのアクセスに困難となる場合がある。

#### 【 0 0 0 5 】

また従来のアプリケーションストレージモデルも、保存されたデータの信頼性を十分に確認することができない場合がある。例えば、従来のオペレーティングシステムは、一般的に初期設定で1つのストレージデバイス上に1つのアプリケーションデータのコピーを保存し、データの冗長性が望まれる場合に、ユーザ又はアプリケーションに自らのアプリケーションデータのコピーを作成し管理することを要求する。個々のストレージデバイス又は第三者のソフトウェアがある程度の冗長性を提供する場合もあるが、アプリケーションが利用可能なストレージリソースがアプリケーション装置によって大いに異なるため、これらの特徴はアプリケーションにとって一貫して利用可能ではない場合がある。またオペレーティングシステム介在の従来のストレージモデルも、プラットフォームを越えるデータへのアクセス性を制限してもよい。例えば、異なるオペレーティングシステムは、異なる、互換性のないフォーマットで同一アプリケーション用のデータを保存する場合があります、1つのプラットフォーム上（例えばオペレーティングシステム及び基礎をなすコンピュータシステムハードウェア）で実行しているアプリケーションのユーザが、異なるプラットフォーム上で実行しているアプリケーションによって保存されたデータにアクセスすることが困難である場合がある。

#### 【 発明の開示 】

#### 【 課題を解決するための手段 】

#### 【 0 0 0 6 】

分散型、ウェブサービス系のストレージシステムの様々な実施態様を開示する。一実施態様に従い、システムは、ウェブサービスプロトコルに従い、データオブジェクトにアク

10

20

30

40

50

セスするためのクライアント要求を受け取るように構成されるウェブサービスインターフェースを含んでいる。所定のデータオブジェクトへのアクセスのための所定のクライアント要求は、所定のデータオブジェクトと対応するキー値を含んでいる。また該システムは、データオブジェクトの複製を保存するように構成される多くのストレージノードを含み、各複製はそれぞれのロケータ値を介しアクセス可能であり、ロケータ値はそれぞれシステム内で固有である。該システムは、データオブジェクトのそれぞれに対するそれぞれのキーマップエントリを保存するように構成されるキーマップインスタンスをさらに含んでもよく、それぞれのキーマップエントリは、所定のデータオブジェクトのそれぞれの保存された複製に対応するキー値及び各ロケータ値を含む。また該システムは、ウェブサービスインターフェースからデータオブジェクトにアクセスするためのクライアント要求を受け取るように構成されるコーディネータをさらに含んでもよい。所定のクライアント要求に応じて、コーディネータはキー値に対応する1つ以上のロケータ値を認識するためにキーマップインスタンスにアクセスするように構成され、特定のロケータ値に対しては、対応する複製を取り出すために対応するストレージノードにアクセスするように設定されてもよい。

10

#### 【0007】

該システムの特定の実装において、該ウェブサービスインターフェースは、ウェブサービスプロトコルに従い、データオブジェクトを保存するためのクライアント要求を受け取るようにさらに設定されてもよく、特定のデータオブジェクトを保存するための特定のクライアント要求は本特定データオブジェクトと対応するキー値を含む。該コーディネータは、ウェブサービスインターフェースからデータオブジェクトを保存するためのクライアント要求を受け取るようにさらに設定されてもよく、特定のクライアント要求に応じて、該コーディネータは1つ以上の対応するストレージノードへ特定のデータオブジェクトの1つ以上の複製を保存するように設定されてもよい。特定のデータオブジェクトの所定の複製を保存することに応じて、所定のストレージノードは、コーディネータに、所定の複製に対応するロケータ値を返すように設定されてもよい。

20

#### 【0008】

本発明は、種々の修正及び代替形態が可能であるが、その特定の実施形態は図面によって一例として示され、本願に詳述される。しかしながら、図面及び詳細な説明は、本発明に開示される特定の形態に制限することを意図せず、逆に、添付の請求項によって定義される本発明の精神及び範囲内に含まれるすべての修正、同等物、及び代替手段を扱うことを理解されたい。

30

#### 【発明を実施するための最良の形態】

#### 【0009】

##### 序文

コンピュータアプリケーションが地理的に分散されるばかりでなく、さらにデータ集中的になるにつれ、アプリケーションデータへの信頼性のある、位置独立型のアクセスの必要性が増加する。例えば、オーサリング、ストレージ、再生アプリケーションなどのマルチメディアアプリケーションは、マルチメディアコンテンツの質と量が向上するにつれ、データストレージ量の増大を必要とする。さらに、データを保存するデバイスの位置に関係ない様々な場所からアプリケーションデータにアクセスすることが望ましい場合がある。例えば、多くのコンピュータがディスク系の相当量のストレージを含むが、これらのストレージに一貫した便利な方法で遠隔からアクセスすることは、技術上及び安全上の困難をもたらす。

40

#### 【0010】

個々のコンピュータを、単独で個々の内部ストレージリソース又はプロビジョニングローカルネットワーク系のストレージリソース（例えば、ネットワーク接続ストレージ（NAS）、ストレージ・エリア・ネットワーク（SAN）など）に依存するように構成されることとは対照的に、インターネット接続データストレージサービスは、例えば、ウェブサービス（WS）プロトコルなどのインターネット系のプロトコルを介し、クライアント

50

に包括的なストレージサービスを提供するように設定されてもよい。ウェブサービスプロトコルなどのインターネット系プロトコルは、一般的に基礎をなすソフトウェア又はハードウェア上で独立して機能するため、一般的にプラットフォーム独立型である。したがって、ウェブサービスとしてデータストレージ機能を提供することは、アプリケーションのホストシステム又はローカルネットワークに実装されるストレージリソースから独立した任意の大きさのストレージへ多くの異なる種類のアプリケーションが直接アクセスすることを可能にする可能性がある。さらに、ウェブサービスアクセス可能ストレージは、一般的にインターネットアクセスを提供する任意の場所からでもアクセス可能であってもよい。ウェブサービスアクセス可能ストレージは、異なるデバイス又はアプリケーションによる共有データへのリモートアクセス、実行中の個々のアプリケーションによる幅広く分散されたデータへのリモートアクセス、共同で作業する分散したユーザ間のデータへのアクセス及びデータの共有、分散したユーザへのアプリケーション結果データの配布、及び多くのその他の同様の機能などの多数の異なるコンピュータ機能の実装を容易にし得る。

#### 【0011】

以下の論考において、ウェブサービス系ストレージシステムにおいて使用されてもよい、可能なデータストレージモデルの一実施形態を説明する。続いて、データストレージモデルに従いストレージサービスを提供するように構成されてもよいストレージサービスシステムを開示し、その様々な要素を詳しく説明する。

#### 【0012】

#### ストレージサービスユーザインターフェース及びストレージモデルの概略

ウェブサービスなどの、サービスとしてのデータストレージをユーザに提供するためのストレージモデルの一実施形態を図1に示す。図示されたモデルにおいて、ストレージサービスインターフェース10は、ストレージサービスに対する対顧客又は対ユーザインターフェースとして提供される。インターフェース10によってユーザに提示されるモデルに従い、ストレージサービスは、インターフェース10を介して任意の数のバケット20a-nとして体系化されてもよい。各バケット20は、任意の数のオブジェクト30a-nを保存するように設定されてもよく、ストレージサービスのユーザによって特定されるデータを入れ替わりに保存してもよい。

#### 【0013】

以下にさらに詳しく述べるように、一部の実施形態において、ストレージサービスインターフェース10は、ウェブサービスモデルに従い、ストレージサービスとそのユーザ間の情報のやり取りをサポートするように設定されてもよい。例えば、一実施形態において、インターフェース10は、例えば`http://storageservice.domain.com`などのUniform Resource Locator (URL)を有するウェブサービスエンドポイントとして、クライアントによってアクセス可能であってもよく、サービスクライアントによって作成されたウェブサービスコールは処理対象となる場合がある。一般的に言えば、ウェブサービスは、Hypertext Transport Protocol (HTTP)又は他の適切なプロトコルのバージョンなどの、1つ以上のインターネット系のアプリケーション層データ転送プロトコルを含む要求インターフェースを介し、要求するクライアントが使用できるように作られた、あらゆる種類のコンピュータサービスを示すことができる。

#### 【0014】

ウェブサービスは、様々な許可サービスプロトコルを使用し、様々なアーキテクチャ形式で実装されてもよい。例えば、Representational State Transfer (REST)形式のウェブサービスアーキテクチャにおいて、ウェブサービスコールに関するパラメータ(例えば、要求されるサービスの種類を特定すること、ユーザ資格、操作すべきユーザデータ、など)は、HTTP GET又はPUT命令など、ウェブサービスエンドポイントに対するウェブサービスコールを呼び出すデータ転送命令に対するパラメータとして特定されてもよい。一部の実装において、各ウェブサービスコールが、外部状態情報を参照することなくそのコールを処理するために必要なすべての情報

を含んでいる場合があるため、R E S T形式ウェブサービスアーキテクチャは処理状態を把握しない。R E S T形式ウェブサービスアーキテクチャとは対照的に、文書系又はメッセージ系ウェブサービスアーキテクチャは、ウェブサービスエンドポイントに送信され、解読され、エンドポイントによって作動し得る文書としてウェブサービスコールに関するパラメータ及びデータをコード化してもよい。例えば、e X t e n s i b l e M a r k u p L a n g u a g e ( X M L ) 又は他の適切なマークアップ言語のバージョンを、ウェブサービス要求文書をフォーマット化するために使用してもよい。一部の実施形態において、要求文書をフォーマット化するために使用されるマークアップ言語を、要求の処理を制御するパラメータの範囲を定めるが、他の実施形態において、マークアップ言語のある特徴は、それ自体（例えば、あるタグ）要求処理の態様を直接制御してもよい。さらに、一部の実施形態において、結果として得られる文書は、例えば、エンドポイントによるウェブサービス要求の処理を容易にするために、S i m p l e O b j e c t A c c e s s P r o t o c o l ( S O A P ) のバージョンなどの他のプロトコル内にカプセル化されてもよい。

10

#### 【 0 0 1 5 】

他のプロトコルが、さらにウェブサービスアーキテクチャの様々な実施形態内で用いられてもよい。例えば、W e b S e r v i c e s D e s c r i p t i o n L a n g u a g e ( W S D L ) のバージョンは、潜在的クライアントにインターフェースで接続する要求を公開するためにウェブサービスエンドポイントによって用いられる。ウェブサービスエンドポイントは、U n i v e r s a l D e s c r i p t i o n , D i s c o v e r y a n d I n t e g r a t i o n ( U D D I ) プロトコルのバージョンなどのディレクトリプロトコルを通じ、ウェブサービスを潜在的クライアントに知らせる場合がある。ウェブサービスインターフェースを介するコンピュータサービスの規定に関係する多くの他の種類のプロトコルが存在する場合があり、あらゆる所定のウェブサービス実装は当該プロトコルのあらゆる適切な組み合わせを使用してもよい。

20

#### 【 0 0 1 6 】

一部の実施形態において、インターフェース 1 0 は、ウェブサービスインターフェースの代わり、又はウェブサービスインターフェースに加えて、ウェブサービスインターフェース以外のインターフェースをサポートしてもよいことに留意されたい。例えば、企業は、異なる種類のインターフェース（例えば、企業のイントラネットに対してカスタマイズされた、独自のインターフェース）を使用するであろう企業内のユーザのみならず、ウェブサービスプロトコルを介しサービスにアクセスするであろう企業外部のクライアントによって使用されるストレージサービスを実装してもよい。一部の実施形態において、インターフェース 1 0 は、それを介してストレージサービスのあらゆるユーザがサービスにアクセスできるようなインターフェースで接続する様々な種類のそれぞれのプロトコルをサポートしてもよい。他の実施形態において、インターフェースの異なる例は、異なるインターフェースアプローチのために提供されてもよい。一部の実施形態において、クライアントとの情報のやり取り（例えば、サービス要求の受け取り及びサービス要求への応答）を取り扱うことに関係するインターフェース 1 0 のこれらの態様は、ストレージサービス（例えば、バケット及びオブジェクトの階層へのサービスの体系）の一般的なアーキテクチャを実装するこれらの態様から独立して実施されてもよいことに留意されたい。このような一部の実施形態において、図 2 の説明とともに、さらに以下に詳しく説明するように、クライアントの情報のやり取り（例えば、ウェブサービスプロトコルを介した）のインターフェース 1 0 の一部は、企業の内部などの特定のユーザによって迂回されてもよい。

30

40

#### 【 0 0 1 7 】

図 1 に示すように、インターフェース 1 0 は、ストレージサービスユーザにバケット 2 0 へのアクセスを提供する。一般的に言えば、バケット 2 0 はストレージサービスのユーザと関連するオブジェクト名前領域のルートとして機能してもよい。例えば、バケット 2 0 はファイルシステムディレクトリ又はフォルダに類似していてもよい。一部の実施形態において、また個々のバケット 2 0 は、ストレージサービスの使用量計算の基準となる場

50



合がある。例えば、ユーザは請求目的のために1つ以上のバケット20と関連付けられる場合があり、ユーザは、これらのバケット20によって設定された名前領域内に階層的に備わるストレージリソース（例えば、ストレージオブジェクト30）の使用に対して請求されてもよい。

#### 【0018】

図で示した実施形態において、バケット20a-nのそれぞれは、各アクセスポリシー23a-nのみならず、関連するメタデータ21a-nを含む。一般的に言えば、メタデータ21が、所定のバケット20の態様又は性質を表現するために使用される。例えば、メタデータ21は、バケット作成日を識別する情報、その作成者の身元情報、バケットがそれに関連する任意のオブジェクト30を有するか否かの情報、その他の適切な情報を含んでいてもよい。一部の実施形態において、メタデータ21は、バケット20に関連するオブジェクト30の合計サイズ、バケット20及び/又はその関連するオブジェクト30に関するユーザのアクセス履歴、バケット20に関連する請求履歴、又は、バケット20の現行又はこれまでの使用に関係する他の適切な情報などの、バケット20の使用特性を示す情報を含んでもよい。一実施形態において、各バケット20は、ストレージサービスによって、ユーザによって、又は自動的に特定される、それぞれの固有の識別子と関連してもよい。固有の識別子は、メタデータ21内、又はバケット20の別個の性質、又はフィールドとして保存されてもよい。一部の実施形態において、所定のバケット20は、明確な参照、ポインタ、又は他の所定のバケット20と関連するオブジェクト30と対応する情報を含まなくてもよい。むしろ、以下に更に詳しく説明するように、オブジェクト30の位置及び選択は、キーマップとして本書で言及する別個のマッピング設備の使用を通じて実行されてもよい。

#### 【0019】

アクセスポリシー23は、バケット20に関連するオブジェクト30へのアクセスを制御するために必要なあらゆる情報を含む。アクセスポリシー23は、バケット20及びその関連するオブジェクト30へアクセスすることが許可された1人又は複数のクライアント、及びどのくらいの役割でアクセスするかを識別する情報を含んでもよい。例えば、アクセスポリシー23は、1人以上のクライアントのために、ユーザ識別子及び/又は認証読み取りを許可されているかをさらに特定する。アクセスポリシー23は、初期設定又はグループ指向規定（例えば、汎用読み取りアクセスを許可するがオブジェクト30への書き込みアクセスを特定のクライアント又はクライアントグループに制限することによって）、又はあらゆる他の望ましいセキュリティモデルをさらに実行してもよい。

#### 【0020】

図で示した実施形態において、所定のバケット20は、1つ以上のオブジェクト30と関連され、各オブジェクトはそれぞれのメタデータ31及びデータ33を含んでもよい。一般的に言えば、オブジェクト30のデータ33はあらゆる一連のビットと対応してもよい。オブジェクト30内に保存されるビットによって表されるデータの種類の、ストレージサービスに透過的であってもよい。すなわち、ビットはテキストデータ、実行可能プログラムコード、音声、ビデオ、又は画像データ、又はあらゆる他の種類のデジタルデータを表し、ストレージサービスは、オブジェクト30を保存し操作することにおいて、これらの様々なデータの種類のうちから必ずしも識別することはない。一部の実施形態において、データ33のサイズは固定上限（例えば1ギガバイト（GB））に限られてもよいが、他の実施形態において、オブジェクト30は、ストレージサービスに使用可能な物理的ストレージリソースのみに従いサイズの拡大を許可してもよい。

#### 【0021】

バケット21に関連するメタデータ21と同様に、メタデータ31は、対応するオブジェクト30についてあらゆる望ましい記述的な情報を保存するように設定されてもよい。例えば、メタデータ31は対応するオブジェクト30が作成された日付及び/又は時間、オブジェクト30のサイズ、オブジェクト30によって保存されたデータ33の種類（例えば、Multi purpose Internet Mail Extensions (MI

10

20

30

40

50

ME) 基準によって識別されるデータの種類)、又はあらゆる他の種類の記述的情報についての情報を含んでもよい。一部の実施形態において、メタデータ31は、アクセスポリシー情報(例えば、オブジェクト30に対して様々なユーザが有してもよいアクセスの種類を示す許可情報)、オブジェクト費用情報(例えば、オブジェクト30と関連する請求割合又は履歴)、又は、オブジェクト30に起因するあらゆる他の適切な情報又は情報の種類の組み合わせのみならず、オブジェクト30と対応するユーザ対話をしめす使用情報又は履歴情報を保存してもよい。一部の実施形態において、クライアントはメタデータ31として保存されるオブジェクトデータと共にメタデータを提供してもよいが、他の事例においては、メタデータ31は、ストレージサービス特徴(例えば、図2に示され、以下に説明されるストレージサービスシステム)を管理するシステムによって作成されるメタデータを含んでもよい。メタデータ31の一部、全部、又はいずれでもないものが、メタデータの種類、クライアントのアクセス権の特定の規定、又は他の適切な要因にしたがって、オブジェクト30へのアクセス権を有するクライアントにアクセス可能であってもよい。

#### 【0022】

一実施形態において、個々のオブジェクト30は、情報の2つの相異なるアイテム、キー又はロケータの、いずれかを使用するストレージサービスシステム内で識別されてもよい。一般的に言えば、キー及びロケータは異なる手段で解釈してもよいが、キー及びロケータは、全体としてのストレージサービスシステムの名前領域の内容内で解釈されてもよい英数字文字列又は他の種類の記号をそれぞれに含んでもよい。一実施形態において、キーは、特定のバケット20(例えば、新規オブジェクトを保存するためのクライアントからの要求に応じて)内に対応するオブジェクト30が作成された時にクライアントによって特定されてもよい。ユーザによってキーが特定されない場合、キーは、ストレージサービスシステムによって新規オブジェクト30に割り当てられる。そのような実施形態において、特定のバケット20のオブジェクト30と関連する各それぞれのキーは、そのバケット20の名前領域内固有であることが求められる場合がある。一般的に言うと、キーは、対応するオブジェクトがストレージサービスシステム内に存在する限り、クライアントが、対応するオブジェクト30にアクセスする際に介する有効な識別子として存続してもよい。

#### 【0023】

所定のバケット20内で、キーは、従来のオペレーティングシステムのファイルシステムに共通のファイルディレクトリ、又はフォルダ名前領域と類似の階層的オブジェクト名前領域を作成するために使用されてもよい。例えば、クライアントは、固有の識別子050739517を有する特定のバケット20へのアクセス権を読み込む、又は書き込むオブジェクトを与えられてもよい。一実施形態において、クライアントは、バケット内でオブジェクトと対応するバケット名前領域内のキーを作成するために、`http://storageservice.domain.com/050739517`へウェブサービスコールを発行してもよい。例えば、クライアントは、「My Documents/Email/message.txt」というキーを使用するこの特定のバケット内でオブジェクト30が作成されることを特定してもよく、このようなオブジェクト30は、アドレス

`http://storageservice.domain.com/050739517/My Documents/Email/message.txt`

へのウェブサービスコールを使用してアクセスされてもよい。

#### 【0024】

一部の実施形態において、キーによって暗示される階層的構造は、オブジェクトストレージの基礎をなす階層に必ずしも反映されなくてもよいことに留意されたい。例えば、一実施形態において、所定のバケット20と関連するオブジェクト30は、オブジェクト30と関連するキーが階層を暗示しているにもかかわらず、ストレージサービスシステム内で平滑な、非階層な形で保存されてもよい。すなわち、このような実施形態において、バケット20は、他のバケット20を階層的に含まなくてもよい。しかしながら、他の実施

形態において、他のバケット 20 内のバケット 20 の階層的包含は、オブジェクトキーによって暗示される階層に対して直接一致する必要はないが、バケットのいかなる当該階層がサポートされてもよい。

#### 【 0 0 2 5 】

一実施形態において、キーによって識別されたオブジェクト 30 へアクセスするためのクライアントによる要求は、要求されたオブジェクト 30 の基礎をなすデータ 33 が取り出される又は修正される前に、クライアント認証手順、アクセス制御チェック、及び / 又はマッピング過程（以下にさらに詳しく説明するような）を受ける場合がある。例えば、クライアントは、クライアントの身元を証明するためにパスワード又は他の資格を提供することを要求される場合があり、一度識別されると、要求されるバケット 20 と関連するアクセス制御パラメータは、識別されたクライアントが要求されたキーへのアクセスを保証する十分な権限が与えられているかどうかを判断するために評価されてもよい。それに反して、ストレージサービスシステムは、キーよりもむしろロケータによってオブジェクト 30 にアクセスする代替の方法をサポートしてもよい。一般的に言うと、ロケータは、ストレージサービスシステムに既知のすべてのオブジェクト 30の中から、オブジェクト 30 の世界的に固有の識別子を表してもよい。すなわち、キーが、特定のバケット 20 と関連する名前領域に対して固有であってもよいが、ロケータはすべてのバケット 20 内のすべてのオブジェクト 30 の世界的な名前領域内で固有であってもよい。例えば、ロケータは、他のロケータの中で固有であるべきストレージサービスシステムによって作成される英数文字列を含んでもよい。さらに以下に詳しく説明するように、一部の実施形態において、オブジェクト 30 の複数インスタンスは、例えば、データ冗長性及びフォルトトレランスを増加するために、ストレージサービスシステムを実装するために使用された物理的ストレージデバイスの全体を通じて複製されてもよい。そのような実施形態において、固有のロケータは、所定のオブジェクト 30 の各複製されたインスタンスに対して存在してもよい。

#### 【 0 0 2 6 】

一部の実施形態において、キーは、オブジェクト 30 がストレージサービスシステム内に存在する限り、オブジェクト 30 へのアクセスに有効でありつづけることが保証されてもよく、当該保証はそのオブジェクト 30 のあらゆる所定のロケータに適応されてもよく、又は適応されなくてもよい。例えば、オブジェクト 30 の複製されたインスタンス（又は複製）が、異なる物理的なストレージの位置（例えば、その基礎となるストレージ媒体の障害又は置換によって）に移動する場合、その新規の位置におけるオブジェクト 30 の移動したインスタンスと対応する他のロケータが作成され、使用されることがあったとしても、特定のインスタンスを参照するロケータは有効ではなくなってもよい。キー及びロケータ間の関係における更なる詳細は、キーマップシステムコンポーネントの操作についての論考で提供される。

#### 【 0 0 2 7 】

キー系対ロケータ系オブジェクトアクセスの例として、オブジェクト 30 は、上記キーによって参照され、[http://storageservice.domain.com/050739517/My Documents/Email/message.txt](http://storageservice.domain.com/050739517/MyDocuments/Email/message.txt) は、ストレージサービスシステム内に保存される 1 つ以上のインスタンスを有することがあり、その 1 つは、<http://storageservice.domain.com/locator/3859C89A208FDB5A> という形式のロケータによって識別されてもよい。

この特定の実施形態において、オブジェクト 30 に対するキー参照は特定のバケット 20 に関連して表現されるが、ロケータ参照は、世界的ロケータ空間（他の種類のロケータコード化又は形式が用いられてもよいが）内の 128 ビットの 16 進絶対数として表現されることに留意されたい。一実施形態において、ロケータに命令されたクライアント発行のウェブサービス要求は、すべての認証、アクセス権、翻訳、又はキー系のウェブサービス要求に適応されてもよい他のステップのいくつか、又はすべてを回避してもよい。処理の

10

20

30

40

50

より少ない層のため、一部のかかる実施形態において、ロケータ系の要求はキー系の要求よりも更に早く処理されてもよい。しかしながら、セキュリティ対策は、ロケータ系要求のために回避されることがあり、クライアントは、敏感なオブジェクト30が改ざんされていないことの個人的保証を（例えば、ロケータの送信及び受信において暗号化された、又は他の方法を使用し）提供する必要がある場合がある。さらに、ロケータの持続は保証されないため（例えば、上記記載のオブジェクトインスタンスの移動の場合）、ロケータ系オブジェクトアクセスの操作を選択するクライアントは、例えば、プリエンプティブ系において、又は存在するロケータがもはや有効ではないことを発見することに応じて新規ロケータを取得することによって、ロケータが使用中に無効になる可能性を許容することが必要となってもよい。

10

#### 【0028】

クライアントのストレージの必要性及び上記の注意点によって、ロケータ系アクセスは、キー系アクセスに関連してより一層の処理能力（例えば、ウェブサービス要求処理の待ち時間及び処理量において）を提供してもよい。例えば、クライアントは、特に敏感でない頻繁にアクセスされるオブジェクト30を参照するためにロケータ系アクセスを使用することを選択してもよい。一部の実施形態において、ロケータ系アクセスは個々のオブジェクト30に基づき不可能であることがあり、したがって、当該のオブジェクトにアクセスすることを望むクライアントに、キー系要求を使用し、当該の要求と関連するあらゆる認証及びアクセス権制御を提出するように強制することに留意されたい。しかしながら、ロケータ系アクセスが可能になったオブジェクト30に対しても、有効なロケータを保持しない悪意のある、又は機能不良のクライアントは、あらゆる所定のオブジェクト30にうまくアクセスする任意の可能性だけを有してもよい。このような可能性は、大規模ロケータ名前領域、ロケータを作成する安全技術（例えば、オブジェクトデータの安全ハッシュの使用）、又は他の適切な技術の使用を通じて任意に防ぐことができる。

20

#### 【0029】

##### ストレージシステムアーキテクチャ及び実装

図1に示すようなウェブサービス系のストレージサービスを実装するように設定されてもよい、ストレージサービスシステムアーキテクチャの一実施形態を図2に示す。図で示した実施形態において、多くのストレージクライアント50a-nは、ネットワーク60を介してウェブサービスプラットフォーム100と情報をやり取りするように設定されてもよい。ウェブサービスプラットフォーム100は、ストレージサービスコーディネータ120（又は単に、コーディネータ120）の1つ以上のインスタンスとインターフェースで接続するように設定されてもよく、ストレージサービスコーディネータは、1つ以上のキーマップインスタンス140及びビットストアノード160とインターフェースで入れ替わりに接続してもよい。さらに、レプリケータ180は、レプリケータキーマップインスタンス190のみならず、ビットストアノード160とインターフェースで接続するように設定されてもよい。コーディネータ120とレプリケータ180の両方は、ノードピッカー130とインターフェースで接続してもよい。図で示した実施形態において、ノードピッカー130、キーマップ140、ビットストアノード160及びレプリケータキーマップ190の各インスタンスは、発見及び障害検出デモン（DFDD）110のそれぞれのインスタンスと関連してもよい。所定のコンポーネントの1つ以上のインスタンスが存在する場合には、以下で言及するコンポーネントは単数又は複数のいずれかであってもよい。しかしながら、いずれの形式の使用も他方を不可能にすることを意図しない。

30

40

#### 【0030】

様々な実施形態において、図2に示すコンポーネントは、コンピュータハードウェア（例えば、マイクロプロセッサ、又はコンピュータシステム）、又はこれらの技術の組み合わせによって指示が直接、又は非直接実行可能である場合、コンピュータハードウェア内で直接実装されてもよい。例えば、図2のコンポーネントは、図29に示すコンピュータシステムの実施形態、及び以下に説明するようなコンピュータ関係のノード（又は単にノード）の多くを含む分散型システムによって実装されてもよい。さまざまな実施形態におい

50

て、所定のストレージサービスシステムコンポーネントの機能性は、特定のノード又はいくつかのノードに渡って分散されることによって実装されてもよい。一部の実施形態において、所定のノードは、1つ以上のストレージサービスシステムコンポーネントの機能性を実装してもよい。図2のコンポーネントの一般的な機能性の概略及び図3に示すようなストレージサービスシステムの典型的な物理的配置に続き、特定のストレージシステムコンポーネントの一部の実施形態を図4から図28に説明とともに以下に提供する。

#### 【0031】

一般的に言えば、ストレージクライアント50は、ネットワーク60を介するウェブサービスプラットフォーム100にウェブサービス要求を提出するように構成されるクライアントのあらゆる種類を網羅する。例えば、所定のストレージクライアント50は、ウェブブラウザの適切なバージョン、又はプラグインモジュール、又はウェブブラウザによって提供される実行可能環境に対する、又は環境内の拡張子を実行するように構成される、他の種類のコードモジュールを含んでもよい。あるいは、ストレージクライアント50は、データベースアプリケーション、メディアアプリケーション、オフィスアプリケーション又は持続的ストレージリソースを使用するあらゆる他のアプリケーションなどのアプリケーションを網羅してもよい。一部の実施形態において、当該アプリケーションは、すべての種類のウェブ系データに対する完全なブラウザサポートを必ずしも実装することのない、ウェブサービス要求を作成し、処理するための十分なプロトコルサポート（例えば、`HyperText Transfer Protocol (HTTP)`の適切なバージョンに対する）を含んでもよい。すなわち、ストレージクライアント50は、ウェブサービスプラットフォーム100と直接対話するように構成されるアプリケーションでもよい。以下に説明するように、ストレージクライアント50は、`Representational State Transfer (REST)`形式ウェブサービスアーキテクチャ、文書系、又はメッセージ系ウェブサービスアーキテクチャ、又は他の適切なウェブサービスアーキテクチャに従いウェブサービス要求を作成するように設定されてもよい。

#### 【0032】

他の実施形態において、ストレージクライアント50は、これらのアプリケーションに透過的である方法における他のアプリケーションに対しウェブサービス系ストレージへのアクセスを提供するように設定されてもよい。例えば、ストレージクライアント50は、上記に記載されるストレージモデルの適切な改良型に従い、ストレージを提供するためにオペレーティングシステム又はファイルシステムを統合するように設定されてもよい。しかしながら、オペレーティングシステム、又はファイルシステムは、ファイル、ディレクトリ、及び/又はフォルダの従来型ファイルシステム階層などの、アプリケーションに対する異なるストレージインターフェースを提示してもよい。このような実施形態において、アプリケーションは、図1のストレージシステムサービスモデルを使用するように修正される必要がなくてもよい。代わりに、ウェブサービスプラットフォーム100へのインターフェース接続の詳細は、オペレーティングシステム環境内で実行されるアプリケーションに代行して、ストレージクライアント50及びオペレーティングシステム、又はファイルシステムによって統合されてもよい。

#### 【0033】

ストレージクライアント50は、ネットワーク60を介するウェブサービスプラットフォーム100へのウェブサービス要求を送り、応答を受信してもよい。様々な実施形態において、ネットワーク60は、クライアント50とプラットフォーム100間のウェブ系の通信を確立するために必要なネットワーキングハードウェア及びプロトコルのあらゆる適切な組み合わせを網羅してもよい。例えば、ネットワーク60は、インターネットを集合的に実装する様々な電気通信網及びサービスプロバイダを概して網羅する。ネットワーク60は、公共又は私的ワイアレスネットワークのみならず、ローカルエリアネットワーク（LAN）又は広域エリアネットワーク（WAN）などの私的ネットワークを含んでもよい。例えば、所定のクライアント50及びウェブサービスプラットフォーム100の両方は、独自の内部ネットワークを有する企業内においてそれぞれ支給されてもよい。そのよ

うな実施形態において、ネットワーク 60 は、インターネットとウェブサービスプラットフォーム 100 間のみならず、所定のクライアント 50 とインターネット間におけるネットワークリンクを確立するために必要な、ハードウェア（例えば、モデム、ルータ、スイッチ、負荷分散装置、プロキシ、サーバなど）及びソフトウェア（例えば、プロトコルスタック、会計ソフトウェア、ファイアウォール/セキュリティソフトウェアなど）を含んでもよい。一部の実施形態において、ストレージクライアント 50 は、公的インターネットよりはむしろ私的ネットワークを使用してウェブサービスプラットフォーム 100 と通信してもよいことに留意されたい。例えば、クライアント 50 は、ストレージサービスシステムとして同一の企業内で支給されてもよい。そのような場合、クライアント 50 は、私的ネットワーク 60 の全体を通じてプラットフォーム 100 に通信してもよい。このような場合、クライアント 50 は、私的ネットワーク 60 全体を通じてプラットフォーム 100 と通信してもよい（例えば、公的にアクセス可能でない、インターネットベースのコミュニケーションプロトコルを使用してもよい LAN、又は WAN）。

10

#### 【0034】

一般的に言うと、ウェブサービスプラットフォーム 100 は、ストレージサービスシステムによって保存されるオブジェクト 30 にアクセスするための要求のような、ウェブサービス要求を受信し、処理するように構成される、1つ以上のエンドポイントを実装するように設定されてもよい。例えば、ウェブサービスプラットフォーム 100 は、前述の例において使用されたエンドポイントである `http://storage.service.domain.com` を実装するように構成されるハードウェア及び/又はソフトウェアを含んでもよく、このようなエンドポイントに命令された HTTP 系ウェブサービス要求は適切に受信され処理される。一実施形態において、ウェブサービスプラットフォーム 100 は、クライアント 50 からウェブサービス要求を受信し、それをコーディネータ 120、又は処理のためのストレージサービスシステムの他のコンポーネントへ転送するように構成されるサーバシステムとして実装されてもよい。他の実施形態において、ウェブサービスプラットフォーム 100 は、大規模ウェブサービス要求プロセス負荷を動的に管理するように構成される、負荷バランシング及び他の要求管理特徴を実装する多くの別個のシステム（例えば、クラスタトポロジーにおいて）として設定されてもよい。

20

#### 【0035】

様々な実施形態において、ウェブサービスプラットフォーム 100 は、上記に詳述されるように、REST 形式、又は、文書系（例えば、SOAP 系）の種類のウェブサービス要求をサポートするように設定されてもよい。1つの特定の実施形態において、プラットフォーム 100 は、ストレージサービスシステムによって管理されるエンティティの様々な操作をサポートする、特定のウェブサービスアプリケーションプログラミングインターフェース（API）を実装するように設定されてもよい。例えば、プラットフォーム 100 によって実装される API は、バケット 20 又はオブジェクト 30 のリスト（フィルタパターン又は基準に従い任意にフィルタされる）、バケット 20 又はオブジェクト 30 のデータ又はメタデータの取り出し、及びバケット 20 又はオブジェクト 30 の作成又は削除を含む、バケット又はオブジェクト上の基本的なクライアント操作をサポートしてもよい。一部の実施形態において、API は、複数のバケット 20 又はオブジェクト 30 に対する操作のバッチアプリケーションなどの、さらに洗練されたクライアント操作をサポートしてもよい。

30

40

#### 【0036】

クライアントのウェブサービス要求のためのアドレス可能なエンドポイントとして機能することに加え、一部の実施形態において、ウェブサービスプラットフォーム 100 は、様々なクライアント管理特徴を実装してもよい。例えば、プラットフォーム 100 は、要求するクライアント 50 の身元を追跡し、クライアント要求の数及び/又は頻度、クライアント 50 に代行して保存された、又は取り出されたオブジェクト 30 のサイズ、クライアント 50 によって使用された全体的なストレージ回線容量、クライアント 50 によって要求されたストレージのクラス、又はあらゆる他の測定可能なクライアント使用パラメー

50

タなどの、ストレージリソースを含む、ウェブサービスのクライアント使用の計測及び会計を統合してもよい。プラットフォーム 100 は、財務会計及び請求システムをさらに実装、又はクライアント使用活動の報告及び請求のための外部システムによって要求又は処理される使用データのデータ系を維持してもよい。

#### 【0037】

一部の実施形態において、プラットフォーム 100 は、クライアント 50 から受信した要求の割合及び種類を反映する評価指数、当該要求によって使用された回線容量、当該要求に対する遅延処理手続きシステム、システムコンポーネント使用（例えば、ストレージサービスシステム内のネットワーク回線容量及び／又はストレージ使用）、要求の結果によりエラーの割合及び種類、要求されたオブジェクト 30 の特徴（例えば、サイズ、データの種類など）、又はあらゆる他の適切な評価指標をなどの様々なストレージサービスシステム操作評価指標を収集及び／又は監視するように設定されてもよい。このような実施形態において、プラットフォーム 100 は、例えば、平均超過時間として、又は様々な分析に従ってもよい特定のデータポイントとして、当該評価指数を総計に収集するように設定されてもよい。様々な実施形態において、当該の評価指数は、クライアント 50 に可視である、又は不可視であってもよい方法でシステムの能力を試験又は監視するために用いられてもよい。例えば、一実施形態において、当該評価指数は、システムコンポーネントを調整、及び維持するためのシステム管理者によって使用されてもよいが、他の実施形態において、当該評価指数（又は当該評価指数と同等の一部）は、ストレージサービスシステムの使用を当該クライアントが監視できるようにクライアント 50 に使用可能にされてもよい。

#### 【0038】

一部の実施形態において、プラットフォーム 100 は、ユーザ認証及びアクセス制御手順をさらに実装してもよい。例えば、所定のバケット 20 と関連する特定のオブジェクト 30 へアクセスするための所定のウェブサービス要求について、プラットフォーム 100 は、要求に関連するクライアント 50 が、所定のバケット 20 及び特定のオブジェクト 30 へのアクセスする権限が与えられているかどうかを確かめるように設定されてもよい。プラットフォーム 100 は、例えば、身元、パスワード、又は所定のバケット 20 に関連する資格者に対する他の資格の評価すること、特定のオブジェクト 30 への許可できる操作を指定する、アクセス制御リストに対する特定のオブジェクト 30 に対し要求されたアクセスを評価することにより、当該権限を判断してもよい。クライアント 50 が、バケット 20 へのアクセスをする、又はオブジェクト 30 に要求された操作（例えば、クライアント 50 が、読み込みアクセス特権のみを有するオブジェクト 30 の書き込みを試みる）を行う十分な資格を有しない場合、プラットフォーム 100 は、例えば、要求するクライアント 50 にエラー状態を示す応答を返すことにより、対応するウェブサービス要求を拒否してもよい。一部の実施形態において、各バケット 20 及びオブジェクト 30 は、そのバケット又はオブジェクトにアクセスすることを統制する、関連するアクセス制御規定を有してもよいことを考慮されたい。当該アクセス制御規定は、メタデータ 21 又は 31 内にアクセス制御情報の記録又はリストとして、あるいはメタデータ 21 及び 31 から別個のデータ構造として保存されてもよい。

#### 【0039】

一部の実施形態において、図 2 のシステムのようなストレージサービスシステムは、任意サイズのオブジェクト 30 をサポートしてもよいが、他の実施形態において、オブジェクト 30 は、チャンクサイズ(chunk size)ともいわれる、ある最大サイズを強制されることもある。このような一部の実施形態において、クライアントが、キーと関連して保存すべきデータを提供し、そのデータがチャンクサイズを上回る場合、プラットフォーム 100 は、チャンクサイズに従い、データを 2 つ以上のチャンクに分配するように設定されてもよい。一実施形態において、プラットフォーム 100 は、関連するキー値を有するそれぞれのオブジェクト 30 として各チャンクを作成するように設定されてもよい。プラットフォーム 100 は、クライアント提供キーを参照するアクセスのための要求が行われた場

合、元のクライアントデータが、チャンクから再構築されることが出来るような方法で、クライアント提供キーの機能として各チャンクに対するキー値を作成してもよい。例えば、プラットフォーム 100 は、クライアントデータから N チャンクを作成するように構成され、クライアント提供キーに対し、N の相異なるパターンを付加することによってこれらのチャンクのための N 対応キーを作成し、N の相異なるパターンは、N チャンクが作成されたのと同じ順序で辞書編集的に順序付けられてもよい。N チャンクのそれぞれは、以下に説明する技術を使用する相異なるオブジェクト 30 として管理され、元のデータは、クライアント提供キーがプレフィックスであるようなキー値を有するすべてのオブジェクト 30 をリストすることによって、またこれらのオブジェクトをリストされた順序で取り出すことによって作成されてもよい。一部の実施形態において、個々のチャンクは、他のチャンクを阻害することなくアクセスされ、修正され、又は取り除かれることがあり、1 つの大きいオブジェクト 30 としてのデータを維持することに関してシステムの能力を改善してもよい。一部の実施形態において、クライアント 50 は、それが提供するデータオブジェクトをチャンクに分割すべきかどうかを指定することを許可されてもよい。

#### 【0040】

図 2 に示すストレージサービスシステムコンポーネントの多くと同様に、他のコンポーネントからのウェブサービスプラットフォーム 100 の機能性を分離することは、ストレージサービスシステムの保守管理及び全体的な拡張性を改善してもよい。例えば、付加的なハードウェア及びソフトウェア供給源は、他のタスクに割り振られた供給源の付加的なウェブサービス処理負荷を独立して維持する為に特に提供されてもよい。さらに、プラットフォーム 100 と関連するあらゆる供給源の障害の影響は、その特定の機能領域の範囲内にとどめられる場合があり、したがって、障害の隔離と解決を容易にする。しかしながら、一部の実施形態において、プラットフォーム 100 の機能性は、他のコンポーネントにおいて統合されてもよいことを考慮されたい。例えば、コーディネータ 120 は、プラットフォーム 100 と関連するタスクを含むように設定されてもよい。

#### 【0041】

ウェブサービスプラットフォーム 100 は、クライアント 50 がストレージサービスシステムの特徴にアクセスすることを通じ、プライマリインターフェースを示すが、当該特徴にたいする唯一のインターフェースを示す必要がないことに留意されたい。例えば、一部の実施形態において、コーディネータ 120 は、ウェブサービスインターフェースから相異してもよい代替 API をサポートするように設定されてもよい。このような代替 API は、例えば、企業の内部のクライアントが迂回ウェブサービスプラットフォーム 100 に対するストレージシステムを提供することを可能にするために使用されてもよい。ある事例において、プラットフォーム 100 の会計及び/又は信任は、管理クライアントなどの内部クライアントに必要がなくてもよい。

#### 【0042】

コーディネータ 120 は、ウェブサービスプラットフォーム 100 及びストレージサービスシステムの他のコンポーネントとの間の行動を協調させるように設定されてもよい。一実施形態において、コーディネータ 120 の主要な義務は、これらのオブジェクト 30 に命令されたウェブサービス要求に応じ、オブジェクト 30 に対するオブジェクトデータ 33 及びメタデータ 31 の読み込み及び書き込み動作の指示を含んでもよい。例えば、以下に詳しく説明するように、オブジェクト読み込みアクセスは、所定のオブジェクト 30 の複製が保存される、ビットストアノード 160 を示すロケータを取り出すために、キーマップインスタンス 140 へのアクセスを行うことに続き、要求されたデータを読み込むために特定のビットストアノード 160 へアクセスすることを行うことを必要としてもよい。同様に、オブジェクト作成及び修正は、必要であれば、作成又は修正されたロケータを反映するために、様々なビットストアノード 160 へオブジェクト 30 の多数の複製を保存し、キーマップインスタンス 140 を更新することを必要としてもよい。一部の実施形態において、コーディネータ 120 は、キーマップインスタンス 140 及びビットストアノード 160 へのこれらの読み込み及び書き込み操作を行うように構成されてもよい。しかし



ながら、一部の実施形態において、コーディネータ 120 は、作成又は修正時点においてオブジェクト 30 の所望の複製の全数を作成することを行わなくてもよいことに留意されたい。さらに以下に詳しく説明するように、一部の実施形態において、オブジェクト 30 に対する書き込み操作は、コーディネータ 120 が、そのオブジェクト 30 の複製のある数の書き込み（たとえば、2 つの複製）を完了した時に終了したと見なされてもよい。さらに、このオブジェクト 30 の複製は、レプリケータ 180 による帯域外又は非同期操作によって完了されてもよい。すなわち、このような実施形態において、オブジェクト作成又は修正操作の帯域内又は同期部分は影響されたオブジェクト 30 の複製の合計所望数より少ない作成物を含んでもよい。コーディネータ 120 は、キーマップインスタンス 140、ビットストアノード 160、及び他のシステムコンポーネントから相異なるコンポーネントとして示されるが、一部の実施形態において、コーディネータ 120 のインスタンスについて、他のシステムコンポーネント（例えば、単独のコンピュータシステムによって実行されるソフトウェアコンポーネントとして）と共に実装されることは可能であるという点に留意されたい。したがって、本書の定義は、コーディネータ 120 がビットストアノード 160、キーマップインスタンス 140、又は別のコンポーネントヘデータを保存し、又はそこからデータを取り出すことに言及してもよいが、一部の実施形態において、このような処理は共有コンピューティングシステム供給源内で起こってもよいことを理解されたい。

#### 【0043】

図 1 に関し、上記に記載されるとおり、一部の実施形態において、ストレージサービスシステムは、バケット系のストレージモデルを含むことがあり、様々なオブジェクト 30 のためのキーが、管理（例えば、会計、請求）、セキュリティ、又はその他の目的のためにバケット 20 内にグループ分けされてもよい。一実施形態において、コーディネータ 120 は、クライアント 50 からのウェブサービス要求に応じ、様々なバケット関連の操作を処理するように設定されてもよい。例えば、コーディネータ 120 は、以下のバケット操作のいくつか、又はすべてを行うように設定されてもよい。

- バケットの作成：バケット 20 のための新しいバケット名前を作成及び保存する。
- 非空バケットの削除：関連するメタデータ 21 を含む所定のバケット 20、及び所定のバケット 20 内のオブジェクト 30 と関連するすべてのキーを削除する。
- 空バケットの削除：所定のバケット 20 と関連するオブジェクト 30 のキーがない場合のみ、所定のバケット 20 及び関連するメタデータ 21 を削除し、それ以外はエラー状態を返す。
- バケットの書き込み：存在するバケット 20 にデータを書き込む（例えば、メタデータ 21）。
- バケットキーのリストアップ：所定のバケット 20（パターン、正規表現、ワイルドカードなどに従い任意に保存され、又はフィルタされた）と関連するオブジェクト 30 のキーをリストする。
- バケットのリストアップ：所定の加入者（例えば、ユーザ又はクライアント 50）と関連するバケットをリストする。

一部の実施形態において、コーディネータ 120 は、コリジョン作成の低い確率を有する適切な乱数アルゴリズムを使用し、新しく作成されたバケット 20 に対する識別子を作成するように設定されてもよい。他の実施形態において、コーディネータ 120 は、例えば、バケット作成に対するクライアント要求における存在するバケット識別子に関し、固有性に対する要求された識別子を調べることによって、クライアント特定バケット識別子をサポートするように設定されてもよい。

#### 【0044】

上記に記載されるように、オブジェクト 30 のインスタンスは、例えば、オブジェクトデータが、あらゆる所定のノード 160 又はそれと関連するインフラストラクチャの障害をしのぐ可能性を増加させるために、所定のノード 160 又は異なるビットストアノード 160 に複製されてもよい。ストレージサービスシステム内のオブジェクトの複製は、以

10

20

30

40

50

下のように、ノードピッカー 130 及びレプリケータ 180 による図示した実施形態で扱われてもよい管理及び最適化のためのいくつかの機会を呈する。

【0045】

コーディネータ 120 がオブジェクト 30 を書き込む要求を受け取る時、コーディネータは、書き込みが完了したことを表す前にノード 160 の所定の数にオブジェクト 30 を相応して書き込んでもよい。しかしながら、書き込まれるべきオブジェクト 30 に対するノード 160 の数と特定の選択は、異なるストレージ規定検討の数によって異なる場合がある。例えば、書き込み操作が完了したことが見なされる前に、オブジェクト 30 の複製のある最低数（例えば、2 つ又は 3 つ）がうまく書き込まれることを要求することは、可能性のある障害を考慮すると、書き込みデータが永続的であるためには賢明である場合がある。しかしながら、複製の最低数を保存するために選ばれたノード 160 が、障害の異なる可能な場所又は領域から分散されることを確実にすることも望ましい場合がある。例えば、同一のデータセンタに位置付けられるノード 160 は、地理的に隔てられたノード 160 よりも同時に（例えば、自然災害、電力障害などの壊滅的な障害によって）機能しなくなる可能性が高い。

10

【0046】

一般的にストレージノード選択論理といわれることもあるノードピッカー 130 は、コーディネータ 120 及びレプリケータ 180 によってアクセス可能なサービスとして構成されることがあり、一実施形態において、様々なストレージ規定を満たすようなオブジェクト読み込み及び書き込み操作のためにノード 160 の選択のためのアルゴリズムを実装してもよい。例えば、上記に概説したオブジェクト 30 の書き込みの場合において、ノードピッカー 130 は書き込みプラン、又はオブジェクト 30 を書き込むべきノード 160 の特定の配列を構築するために操作してもよい。特定の書き込みプランを立てることに  
いて、ノードピッカー 130 は、書き込みプランが成功の妥当な可能性を有することを確実にするように設定されてもよく、例えば、書き込みプランで特定されたノード 160 は、実際に操作可能であり、オブジェクト 30 を受け入れるために使用可能である十分なストレージリソースを有することが期待され、完了した場合、その書き込みプランは、書き込み操作に関連するすべてのストレージ規定を満たす。書き込みストレージ規定の例は以下を含む。

20

- 永続性規定：書き込みプランがうまく完了した場合、オブジェクト 30 のインスタンスは、少なくとも N 相違ノード 160 に保存される。
- 領域多様性規定：可能であれば、書き込みプランは、少なくとも M 相違領域間に分散されるノード 160 を含む。
- ローカル性規定：可能であれば、書き込みプランは、要求するコーディネータ 120 にたいしてローカルである領域のノード 160 に、選択（例えば、数で）を与える。
- 負荷平衡規定：ノード 160 間の書き込み要求トラフィックを均一にするよう試みる（例えば、「ホットノード」を避けるために）。
- 空間平衡規定：ノード 160 間の保存供給源容量の使用を均一にするよう試みる。
- 最低コスト連鎖規定：書き込みプランにおけるノード書き込み操作の配列の合計コスト（例えば、ネットワーク遅延）を最小限にするよう試みる。

30

40

【0047】

様々な実施形態において、ノードピッカー 130 は、所定の書き込みプランを立てる時、これらの規定のいくつか又はすべて、又はリストアップされない他の規定を考慮に入れることに留意されたい。さらに、異なる規定が、異なる優先事項によって重要性を増す場合がある。例えば、一実施形態において、持続性規定は、他の規定が最善努力系で満たされることがある一方、すべての書き込みプランが満たされなければならない義務的規定であってもよい。一部の事例において、一部のストレージ規定は他と拮抗してもよい。例えば、異なる領域間にオブジェクトインスタンスの幅広い分散を助ける領域多様性規定は、概して特定の領域内にオブジェクトインスタンスをローカル化するローカル性規定に反する。オブジェクトインスタンスの数が十分に大きい場合、両規定を満たすことが可能であ

50

る場合がある。例えば、オブジェクト 30 の 5 つのインスタンスが作成された場合、要求するコーディネータ 120 にローカル的である 2 つのインスタンスを 2 つの相異なる領域に保存し、3 つのインスタンスを 3 番目の相異なる領域に保存し、したがって、ローカル性規定及び領域多様性規定の両方を満たすことが可能である場合がある。書き込みプランを特定するすべての規定を満たすことが不可能でなければ、ノードピッカー 130 は、満たされるこれらの規定を優先させるよう試み、最善書き込みプランを作成し、又はオブジェクト書き込みが十分に行われなことを示す要求するコーディネータ 120 に対するエラー表示を返してもよい。

#### 【0048】

また一部の実施形態において、ノードピッカー 130 は、読み込みオブジェクト 30 のコーディネータ 120 を援助してもよい。例えば、オブジェクト読み込み操作は、もともと、又は最も最近に書いた要求されたオブジェクト 30 で一部のコーディネータ以外のコーディネータ 120 によって要求されてもよい。したがって、書き込んでいるコーディネータ 120 に関連しローカルで保存されてもよいオブジェクト 30 のインスタンスは、読み込んでいるコーディネータ 120 に関連してローカルではなくてもよい。ノードピッカー 130 は、読み込んでいるコーディネータ 120 に使用可能な読み込み操作を提供し得るノード 160 を識別するように設定されてもよい。例えば、ノードピッカー 130 は、読み込んでいるコーディネータ 120 に最も近い（例えば、地理的距離又はネットワークポロジに関して）ノード 160、又は最も高い読み込み回線容量を提供するノード 160（例えば最小負荷ノード 160 又はストレージハードウェアのより高い能力クラスを有するノード 160）を識別してもよく、又はノードピッカー 130 は、オブジェクト 30 を読み込むノード 160 を選択するための他の性能基準を使用してもよい。他の実施形態において、コーディネータ 120 の読み込みに関する読み込みの操作を最適化するよりは、ノードピッカー 130 は、全体としての（例えば、広域読み込みスループットを最大限にするために）性能を最適化するために、同時読み込み操作を広域的にプランしてもよい。

#### 【0049】

書き込みプランを立て、オブジェクト読み込み操作に関連してコーディネータ 120 に通知するために、ノードピッカー 130 は、例えば、操作状態及び利用可能な供給源に関して、ノード 160 の状態を監視するように設定されてもよい。一実施形態において、ノードピッカー 130 は、現在操作可能なストレージサービスシステム内のノード 160 を識別するために、DFDD110（以下に説明する）のインスタンスと交流しあうように構成されてもよい。ノードピッカー 130 が操作可能ノード 160 に気付くと、これらのノードに、それぞれが利用可能な供給源（例えば、ストレージ容量）を確実にするためにクエリを行ってもよい。ノード 160 の操作状態及び供給源状態は時間と共に変化してもよい。一部の実施形態において、ノードピッカー 130 はしばしば DFDD110 を介し、操作状態の情報を新しくし、これに伴うノード 160 がこの供給源状態情報を新しくすることを集計してもよい。一部の例において、ノードピッカー 130 は、ノード 160 の状態の完全な同期視野を有しない場合があることに留意されたい。例えば、ノードピッカー 130 によって使用可能であると思われる特定のノード 160 は、実際には状態情報の最近の更新のために機能しない場合がある。このような例において、ノードピッカー 130 は、コーディネータ 120 によって完了されることが可能である場合がある読み込み、又は書き込みプランを保証することが不可能である場合がある。コーディネータ 120 が、ノードピッカー 130 によって特定されるノード 160 へアクセスできない場合、関連の操作は機能しない場合があり、そっくりそのままコーディネータ 120 によって再度試みられるか、又はコーディネータ 120 は、ノードピッカー 130 が要求されたプランを改訂することを交渉してもよい。一部の事例において、義務的規定を満たすことが以前可能であるが、書き込みプランにおいて特定されるノード 160 の障害が、任意の又は最善努力のストレージ規定のみに影響を及ぼす場合、書き込みプランを完了することが可能であってもよい。このような一部の実施形態において、レプリケータ 180 は、以下に

10

20

30

40

50

説明するように、暫く経ってから満たされなかったストレージ規定を満たすための試みをするように設定されてもよい。

【 0 0 5 0 】

一部の実施形態において、ノードピッカー 1 3 0 の複数のインスタンスは、ストレージサービスシステムにわたって分散されてもよい。例えば、ノードピッカー 1 3 0 のそれぞれのインスタンスは、コーディネータ 1 2 0 の各インスタンス用に分散してもよい。ノードピッカー 1 3 0 は、API を介し、コーディネータ 1 2 0 ( 及び、レプリケータ 1 8 0 ) からアクセスされてもよいサービスとして分散されてもよく、この構造は必須ではない。他の実施形態において、ノードピッカー 1 3 0 の機能性は、コーディネータ 1 2 0 及び / 又はレプリケータ 1 8 0 のインスタンス内に直接合体されてもよい。

10

【 0 0 5 1 】

上記に記載されるように、オブジェクトデータの信頼性及び利用可能性は、ストレージデバイスシステムにわたるオブジェクト 3 0 を複製することによって増加してもよい。例えば、地理的に分散されたシステム内でオブジェクト 3 0 のインスタンス又は複製を分散することは、当該クライアントに近い一部のオブジェクトインスタンスを利用可能に配置することによって、当該オブジェクト 3 0 へアクセスを試みる同様に分散されたクライアント 5 0 の機能を改善してもよい。( オブジェクトの複製の内容において、「インスタンス」及び「複製」の用語は本書において同義的に使用される場合があることに留意されたい。 ) さらに、オブジェクト複製は、特定のオブジェクトインスタンスの破壊に起因するデータ損失の可能性を概して減少してもよい。しかしながら、所定の時における一部の実施形態においては、オブジェクト 3 0 の有効な複製の数は、複製の所望の数又は目標とする数よりも少なくてもよい。例えば、ストレージサービスシステムに全体にわたり複製ストレージ規定を実行することは、各オブジェクト 3 0 の複製の特定の目標数 ( 例えば、又は他の適切な数 ) がいかなる所定の時においても存在することを特定してもよい。しかしながら、所定のオブジェクト 3 0 について、有効な複製の実際数は、様々な理由で目標の数より少ない場合がある。例えば、これまで有効な複製は、それが保存されたデバイスの障害のため、アクセス不可能になる場合がある。代わりに、一部の実施形態において、コーディネータ 1 2 0 によって書き込まれたオブジェクト 3 0 のインスタンスの数は、そのオブジェクト 3 0 のための複製の目標の数よりも少なくてもよい。例えば、上記に記載されるように、インスタンスは、ノードピッカー 1 3 0 によって特定される書き込みプランに従い書き込まれてもよく、目標の数よりも少ないインスタンスを要求する持続性規定を考慮に入れてもよい。

20

30

【 0 0 5 2 】

一実施形態において、レプリケータ 1 8 0 は、各オブジェクト 3 0 の有効な複製の数が目標の数を満たすかどうか ( 例えば、判断が下される時点において、複製の数が、少なくとも目標の数であるかどうか ) 決定するために、オブジェクト 3 0 を調べるために操作されてもよい。具体的には、一実施形態において、レプリケータ 1 8 0 は、各オブジェクト 3 0 のインスタンスの数及び場所を特定する記録に渡り、継続して反復するように設定されてもよい。例えば、レプリケータ 1 8 0 は、以下に詳しく説明するキーマップインスタンス 1 4 0 のように、オブジェクトキーと複製されたオブジェクトインスタンスを識別する対応するロケータ間のマッピングを保存するように設定されてもよい、レプリケータキーマップ 1 9 0 を参照してもよい。( 他の実施形態において、レプリケータ 1 8 0 は、キーマップの専用インスタンスよりは、1 つのキーマップインスタンス 1 4 0 を参照してもよい。 ) 一部の実施形態において、複数のレプリケータ 1 8 0 は、キーマップ空間の異なる部分を同時に調査するように設定されてもよく、ストレージサービスシステムによって管理されるすべてのオブジェクト 3 0 の状態を調べるために必要とされる全体の時間を削減してもよい。

40

【 0 0 5 3 】

レプリケータ 1 8 0 が、所定のオブジェクト 3 0 に対する有効な複製の目標とする数が満たされないと判断する場合、所定のオブジェクト 3 0 への書き込み操作を行うコーディ

50

ネータ 1 2 0 と同様の方法で、レプリケータは所定のオブジェクト 3 0 の追加的な複製を書き込むように設定されてもよい。例えば、レプリケータ 1 8 0 は、上記に記載されるように、付加的レプリケータの作成のための書き込みプランを取得するためにノードピッカー 1 3 0 とインターフェースで接続してもよい。代わりに、レプリケータ 1 8 0 は、オブジェクト複製を作成するための規定を反映する、独自のアルゴリズムを実行してもよい。一部の実施形態において、レプリケータ 1 8 0 は、付加的な複製が要求される状態下によって、オブジェクト 3 0 のための複製を作成する異なる優先事項と一致してもよい。例えば、レプリケータキーマップ 1 9 0 にリストされたロケータの目標の数よりも少ない数を有するオブジェクト 3 0 は、コーディネータ 1 2 0 によって最近書き込まれた場合がある。反対に、いくつかが無効であるロケータの目標とする数を有するオブジェクト 3 0 は、根本的なストレージの障害を提示している場合がある。規定の問題として、レプリケータ 1 8 0 は、後者の事例の前に前者の事例を、訂正することを試みてもよく、その逆もまた同様である。代わりに、レプリケータ 1 8 0 は、該状態を生じさせる特定の状況に関わらず、この状態に遭遇する場合はいつでも、有効な複製の目標とする数よりも少ない数を有するあらゆるオブジェクト 3 0 のための付加的な複製を作成することを試みてもよい。

#### 【 0 0 5 4 】

上記に記載されるように、オブジェクト 3 0 の全体的な信頼性は、例えば、異なる領域又はデータセンタ内に、オブジェクトデータの複製を保存することによって増加してもよい。しかしながら、一部の実施形態において、各複製は、オブジェクトデータの正確なコピーと対応する必要があることに留意されたい。一実施形態において、オブジェクト 3 0 は、オブジェクトデータが、すべての作成された部分を使用しなくても再度作成できる、冗長するコード化スキーム（パリティ、エラー訂正又は他のスキームなど）に従い多くの部分又は「破片」に分割されてもよい。例えば、オブジェクト 3 0 から N 部分を作成するために様々なスキームを使用することで、オブジェクトデータは、コード化スキームに従い、部分の任意の N - 1、N 部分の任意の過半数、又は部分の他の組み合わせから再度作成されてもよい。このような実施形態において、オブジェクト 3 0 の複製は、作成された部分、又は部分の一部の組み合わせと対応してもよい。このような方法は、オブジェクトデータの複数の完成したコピーを保存することに比べ、データストレージ要求を減少すると同時に、効果的なフォルトトレランスを提供してもよい。しかしながら、一部の実施形態において、冗長するコード化技術は、オブジェクトデータの完成した複製との組み合わせにおいても使用される場合があることに留意されたい。例えば、オブジェクトデータの複数の個々の完成したコピーは、上記に記載されるように、適切な冗長のコード化技術に従い、決定された複数部分のそれぞれの収集物としてノード 1 6 0 間に保存されてもよい。最後に、一部の実施形態において、一部のオブジェクト 3 0 は、複製又はフォルトトレランスのいかなる度合いをも有して保存される必要があることに留意されたい。例えば、ストレージクラスの説明とともに以下に説明するように、クライアントは、フォルトトレランスの高い程度を特定するストレージクラスについておそらく低コストである、ある程度の又は全く程度を有しないフォルトトレランスを特定するストレージクラスに従いオブジェクト 3 0 が保存されることを要求してもよい。

#### 【 0 0 5 5 】

一般的に言えば、キーマップインスタンス 1 4 0 は、オブジェクト 3 0 のキーと特定のインスタンスのロケータ、又はオブジェクト 3 0 の複製間の関係の記録を提供してもよい。このような記録を保存することにおいては、キーマップインスタンス 1 4 0 は、オブジェクト 3 0 がストレージシステム内に複製された度合いも反映する（例えば、いくつかのオブジェクト 3 0 のインスタンスが存在するか、どのようにそれらが参照される場合があるか）。ビットストアノード 1 6 0 は、概して、ロケータによって識別されるオブジェクト 3 0 の個々のインスタンスのためのストレージを提供してもよい。しかしながら、所定のノード 1 6 0 は、あらゆる他のノード 1 6 0 に関するインスタンスの状態、又は対応するオブジェクト 3 0 のインスタンスのロケータ及びキーとの間の関係に気付かないことがある。すなわち、一般的に言うと、キーマップインスタンス 1 4 0 によって維持される状態

情報は、ビットストアノード 160 に透過的であることがある。DFDD110 は、コーディネータ 120 及びレプリケータ 180 などの DFDD110 のクライアントが正確であるが、遅延の可能性が一部の検出された状態のビューを得るために、ノード 160 及び / 又はキーマップインスタンス 140 ( 及び、実行されている場合、レプリケータキーマップ 190 ) の操作状態についての状態情報を検出し通信する操作をしてもよい。これらのコンポーネントは、以下にさらに詳しく取り扱う。

#### 【0056】

図 2 のストレージサービスシステムアーキテクチャの一部のコンポーネントの物理的展開を示す一実施形態を、図 3 に示す。図示した実施形態において、データセンタ 300 は、2 つの領域である 310 a - b を含んで示される。さらに、領域 310 c - d がデータセンタ 300 の外部から示され、領域 310 a - d は、ネットワーク 60 を介して相互接続されている。領域 310 a - d のそれぞれは、それぞれコーディネータインスタンス 120 a - d を含む。さらに領域 310 a - d は、図 3 に示されない、図 2 の他のコンポーネントのみならず、様々な組み合わせのビットストアノード 160 及びキーマップインスタンス 140 を含んでもよい。例えば、領域 310 a は 4 つのビットストアノード 160 を含み、領域 310 b は 3 つのビットストアノード 160 及びキーマップインスタンス 140 を含み、領域 310 c は 2 つのビットストアノード 160 を含み、領域 310 d は 1 つのビットストアノード 160 及び 1 つのキーマップインスタンス 140 を含む。

#### 【0057】

上記に記載されるように、一実施形態において、領域 310 a - d のそれぞれは、障害から独立している、又はほとんど相関しない場所と考えられてもよい。すなわち、障害が起こっているいかなる所定の領域 310 の可能性は、いかなる他の所定の領域 310 の障害の可能性から独立している、又は相関しないか、又は障害の可能性の相関性はしきい値量よりも低い場合がある。例えば、2 つの領域 310 は、同時に障害をきたす可能性が 10 % 以下を示す。障害の相関関係又は独立性は、あらゆる適切な統計的又は確率的な技術を使用して測定され、様々な方法で実行されてもよい。例えば、領域 310 は、物理的に分割されるか、又は独立したユーティリティグリッドへ接続されてもよく、これはおそらく 1 つの領域 310 に影響をおよぼす大災害が、他に影響を及ぼさないようにする。同様に、データセンタ 300 内で、相異なる領域 310 は、1 つの領域 310 が別の領域 310 の障害に関わらず、操作を継続することを可能にするために機能してもよい、独立したバックアップ電力源、ネットワーク接続、又は他の冗長性供給源を有してもよい。

#### 【0058】

一実施形態において、それぞれの障害の可能性間における、小さいがゼロではない相関間関係を有する 2 つの領域 310 は、障害の独立した可能性を有するとして依然と言われることもあることに留意されたい。例えば、バックアップ電力源、冷却などのためのしっかりとした、独立したシステムをそれぞれ有するにもかかわらず、所定のデータセンタ 300 内の 2 つの領域 310 は、非常に大規模の大災害 ( 例えば、データセンタ 300 の全体を破壊するに十分な爆発 ) の場合における同時障害の影響を受けやすい場合がある。しかしながら、これらの 2 つの領域 310 が同時に機能しなくなるのに十分である事象の可能性は、非常に少なく、実用的な目的として、該 2 つの領域 310 は、障害の独立した可能性を有すると言える場合がある。

#### 【0059】

領域 310 は、階層の付加的レベルを含んでもよい ( 本書に示されない ) 。例えば、一実施形態において、あらゆる適切な領域体系が用いられる場合があるが、領域 310 は、ビットストアノード 160 などの、個々のノードにさらに細分化され、ラックに細分化されてもよい。一般的に言えば、領域 310 は、該領域内で展開されるストレージサービスコンポーネントを実装するために十分なコンピューティング供給源を含んでもよい。例えば、各ビットストアノード 160 は、図 4 から図 9 の説明とともに以下に説明するような様々なハードウェア及びソフトウェアコンポーネントを含む場合が一部の独立したコンピュータシステムとして実装されてもよい。同様に、各キーマップインスタンス 140 は、

図 10 から図 22 の説明とともに以下に説明するような構成の多くのコンピュータシステムを介して実装されてもよい。

#### 【0060】

一部の実施形態において、ウェブサービスプラットフォーム 100、コーディネータ 120、ノードピッカー 130、レプリケータ 180、及び DFDD 110 などのコンポーネントは、コンポーネントが展開される各領域 310 内の個別のコンピューティング供給源を介して実装されてもよい。例えば、これらのコンポーネントのそれぞれは、それぞれの計算システムによって実行可能なインストラクション及びデータの 1 組として実装されてもよい。代わりに、これらのコンポーネントのいくつか、又はすべては、1 つ以上の計算システムで同時に実行する場合が一部の過程として実装されてもよい。一部の実施形態において、これらのコンポーネントのいくつか、又はすべてを実装するために使用されるコンピューティング供給源は、ビットストアノード 160 又はキーマップインスタンス 140 を実装するために使用されるこれらの供給源を共有してもよい。例えば、計算システムは、コーディネータ 120 機能性のみならず、キーマップ 140 機能性の一部分の、両方を実装するように設定されてもよい。一般的に言えば、個々の領域 310 内で展開されるコンピューティング供給源にわたる図 2 のコンポーネントのあらゆる適切なセグメント化が、用いられる場合がある。図 3 に示すように、異なる領域 310 は、ストレージシステムコンポーネントの異なる組み合わせを含む場合があり、示された実施形態は、制限することよりもむしろ実例となることを意図したものであることに留意されたい。

#### 【0061】

さらに、異なるストレージサービスシステムコンポーネントは、通信プロトコルのあらゆる適切な種類に従い、通信してもよい。例えば、図 2 の一部のコンポーネントが個別のアプリケーション又は実行可能な過程として実装される時、これらはオペレーティングシステム、又はプラットフォーム（例えば、遠隔手続きコール、キュー、メールボックス、ソケットなど）によって、提供される場合がある標準のプロセス間の通信を互いに使用して、又は標準の、又は独自のプラットフォーム独立型通信プロトコルの使用によって通信してもよい。当該プロトコルは、ハンドシェイク／肯定応答、エラー検出及び訂正、又は通信コンポーネントに要求される、又は望まれる場合があるその他の通信特徴の任意のレベルをサポートしてもよい、処理状態を把握する、又は処理状態を把握しないプロトコルを含んでもよい。例えば、1 つのストレージサービスシステムの実施形態において、コンポーネント間通信の実質的程度は、Transmission Control Protocol (TCP)、User Datagram Protocol (UDP)、又は同様の標準の、又は独自の輸送プロトコルのバージョンなどの、適切なインターネット輸送層プロトコルを使用して実装されてもよい。しかしながら、ストレージサービスシステムコンポーネント間の通信は、プロトコル抽象化のより高い層にプロトコルを使用して実装される場合があることも考慮にされたい。例えば、クライアント 50 とウェブサービスインターフェース 100 間の通信のように、ストレージサービスシステムコンポーネント間の通信は、例えば、HTTP 上のウェブサービスコールのような、アプリケーション層プロトコルを使用して実行されてもよい。

#### 【0062】

##### ビットストア構造

上記に論じるように、図 2 に示すストレージサービスシステムアーキテクチャにおいて、ビットストアノード 160 は、概してストレージサービスシステムによって管理されるオブジェクト 30 に対するストレージを提供するために動作してもよい。ビットストアノード 160 の 1 つの例示的な実施形態を図 4 に示す。図で示された実施形態において、ビットストアノード 160 は、ストレージリパッカー 163 及び論理ファイルインプット／アウトプット (I/O) マネージャー 165 とインターフェースで接続するように構成されるストレージノード管理 (SNM) コントローラ 161 を含む。マネージャー 165 は、ファイルシステム 167 とインターフェースで接続するように構成され、入れ替わりに 1 つ以上のストレージデバイス 169 を管理するように構成される。様々な実施形態にお

いて、コントローラ 161、ストレージレパッカー 163、論理ファイル I/O マネージャー 165、又はファイルシステム 167 のいずれかは、コンピュータ接続可能媒体に保存され、以下に説明する機能を実行するためにコンピュータによって実行可能であってもよいインストラクションとして実装されてもよい。代わりに、これらのコンポーネントのいずれかは、専用のハードウェアサーキット、又はデバイスによって実装されてもよい。

#### 【0063】

一実施形態において、SNM コントローラ 161 は、API に従い、活動を実現するためにノード 160 の他のコンポーネントの活動を協調するのみならず、ノード 160 のクライアントに対するオブジェクトストレージ API を提供するように設定されてもよい。例えば、コントローラ 120 は、SNM コントローラ 161 によってもたらされる API を介し、所定のノード 160 へのオブジェクト 30 の保存する、及び所定のノード 160 からのオブジェクト 30 の取り出すように設定されてもよい。API 管理が本願で SNM コントローラ 161 の特徴として説明されるが、一部の実施形態において、ノード 160 の API 処理機能はモジュールに、あるいは、SNM コントローラ 161 から相異なるコンポーネントに実装される場合があることを考慮されたい。

#### 【0064】

オブジェクトストレージ API は、オブジェクトを格納、取得、及び解放する操作をサポートしてもよい。1 つのこのような実施形態において、総称的に保存操作、又は書き込み操作といわれることもある、オブジェクト格納操作は、操作の引数又はパラメータとしてオブジェクト 30 のデータ及び/又はメタデータを特定してもよい。所定のノード 160 が完了次第、以下にさらに詳しく説明するように、ストレージサービスシステム全体に渡って保存される他のオブジェクト 30 のすべてに対して取り出す所定のノード 160 のオブジェクトインスタンスを個別に識別してもよい、格納操作は、保存されたオブジェクト 30 と対応するクライアントロケータを要求するクライアントに返してもよい。

#### 【0065】

反対に、総称的に読み込み、又は取り出し操作と言われることもある、オブジェクト取得操作は、パラメータとして、オブジェクト 30 のロケータを特定してもよい。完了次第、取得操作は、特定されたロケータと対応するオブジェクトデータ及び/又はメタデータを、要求するクライアントに返してもよい。一部の実施形態において、該取得操作は、オブジェクトデータ、メタデータ、又はその両方がクライアントに対して返されるべきかどうかを要求するクライアントが特定することが出来るパラメータをサポートしてもよい。

#### 【0066】

取得操作のように、総称的に削除、又は除去操作と言われることもある、オブジェクト解放操作は、パラメータとしてオブジェクト 30 のロケータを特定してもよい。しかしながら、完了次第、解放操作は、参照されるオブジェクト 30 とあらかじめ関連してストレージリソースを解放してもよく、当該供給源は他のオブジェクト 30 を保存するために使用されてもよい。一実施形態において、ロケータが解放されると、ロケータに対するそれに続く取得操作は一定期間続くか、又は続かなくてもよい。すなわち、解放操作は、再使用のためのストレージリソースを解放する場合が一部のノード 160 に対する信号としての役割を示す場合があるが、ノード 160 がそうすることを直ちに試みない、また通知を試みない場合には、クライアントとの当該の再使用を同期化してもよい。したがって、解放に続くオブジェクト 30 へのアクセスのためのクライアントによる持続的な試みは、任意の期間にわたり継続される場合があり、それに続いてオブジェクト 30 は予告無しにアクセス不可能になる場合がある。他の実施形態において、ノード 160 はオブジェクトデータが依然として使用可能であるか否かに関わらず、あらかじめ解放されたロケータへのクライアントのアクセスを妨げるように設定されてもよい。

#### 【0067】

様々な実施形態において、取得及び解放操作は、あらゆる適切なプロトコルに従い、他のパラメータを使用する及び/又は様々な状態、エラー又は他の表示を返すことがあることに考慮されたい。例えば、格納操作は、保存されるべき要求されたオブジェクト 30 に



対するノード160に不十分な供給源がある場合、あるいは、何らかの他の理由で格納が完了されない場合、エラー状態を返すことがある。一部の実施形態において、ノード160のオブジェクトストレージAPIは、他の操作を含んでもよいことを考慮されたい。例えば、APIは、複製操作をサポートすることによってオブジェクト複製の作成を容易にするように設定されてもよい。一実施形態において、目的とするノード160に保存すべきオブジェクト30のデータを供給する代わりに、要求するクライアントが異なるノード160にそのオブジェクト30のロケータを特定することを除き、複製操作は、格納操作と同様に行われてもよい。目標とするノード160は、オブジェクトデータ及び/又はメタデータを得るために、特定されたノード160と交流し、目標とするノードと関連するオブジェクトのロケータをクライアントに返してもよい。他の実施形態において、ノード160は、オブジェクト30での他の適切な操作をサポートしてもよい。

10

#### 【0068】

上記に記載されるように、格納、取得、及び解放操作を実装する、一部の実施形態において、存在するオブジェクト30は同一の場所で改訂されない場合があることに留意されたい。むしろ、オブジェクト30のインスタンスは、修正されたデータを含む新規インスタンスを書き込んだ後に、存在するインスタンスを解放することによって効果的に修正されてもよい。このような手段は、例えば、オブジェクト30への修正のレンダリングが元のサイズよりも小さい、又は大きいことが生じた場合、フラグメンテーションの削減又はオブジェクト再配置によって、ノード160の根本的な管理層の実装を簡素化してもよい。ウェブサービスプラットフォーム100に関して以下にさらに詳しく説明するように、一部の実施形態において、ストレージサービスシステムは、ラージオブジェクトをチャンクに分割し、各チャンクは相異なるオブジェクト30として管理されてもよい。この手段は、再書き込みが必要なことがあるチャンクの範囲を制限することによって頻繁に修正されてもよいラージオブジェクトを処理することにおいて、ノード160の機能を改善してもよい。しかしながら、他の実施形態において、ノード160は、今説明した再書き込みの解放手段を通じるよりはむしろ、所定の位置に、オブジェクト30の修正をサポートするために必要なこれらの特徴を含んでもよいことを考慮されたい。

20

#### 【0069】

図示した実施形態において、論理ファイルI/Oマネージャー165（又は、単にマネージャー165）は、オブジェクト30が備わっている場合がある1つ以上の論理的連続ストレージ空間を、SNMコントローラ161及びリパッカー163にもたらしするために、根本的なデバイス、又はファイルシステム特徴に仮想化するように設定されてもよい。例えば、所定のオブジェクト30は、オフセットからのストレージ空間内及び範囲内のオフセット（例えば、データ及びメタデータを含むオブジェクトサイズに関して）に従い、論理ストレージ空間内に位置されてもよい。そのような論理ストレージ空間を提供することによって、マネージャー165は、当該根本的ストレージの実装の詳細に関わらず、SNMコントローラ161に対する根本的なストレージの均一の視野をもたらし場合がある。

30

#### 【0070】

論理ストレージ空間内でオブジェクト30へアクセスすることを容易にするために、一実施形態において、マネージャー165は、各オブジェクト30をノード160に保存するためのオブジェクトインデックス値（オブジェクトインデックスとしても言及される）を指定するように設定されてもよい。一般的に言えば、あらゆる所定のオブジェクト30のインデックスは、所定のノード160内で固有であってもよい。例えば、一実施形態において、オブジェクトインデックスは、オブジェクト30がノード160に保存される場合には常に、カウンタを増加し、オブジェクトインデックスとしての結果として得られたカウンタ値によって得られてもよい。（複数のオブジェクト書き込み操作が同時に進行することができる実施形態において、カウンタ増加は、例えば、シリアライゼーションを通じて、オブジェクトインデックス値が、一定の予想可能な形に指定されることを確実にするために、同期化されてもよい。）64ビットの符号なし整数などの、十分に大きいカウンタ値は、例えば、実用的な目的のために、各オブジェクト30が固有のインデックス値

40

50

を指定されることを確実にしてもよい。このようなカウンタは、例えば、 $2^{64}$ オブジェクトが保存された後一周してもよく、その後前もって作成されたインデックス値を繰り返してもよい。しかしながら、カウンタが一周した後にノード160内に所定のインデックス値を前もって指定されたオブジェクト30が依然として存在するということは極めて可能性が低いように、コリジョンの可能性は非常に低い。オブジェクトインデックスを指定するためのあらゆる他の適切な方法も用いられる場合があることに留意されたい。以下に説明するように、オブジェクトインデックス値は、特定のオブジェクト30を参照するためのコーディネータ120、又はノード160の他のクライアントによって使用されるロケータ値を決定するために、ノード160の固有の識別子との組み合わせで使用されてもよい。

10

#### 【0071】

マネージャー165は、オブジェクト30が、オブジェクトアクセスを容易にする方法で論理ストレージ空間内に位置される場所についての情報を体系化するために、上記に記載される固有のオブジェクトインデックス値を使用するように設定されてもよい。例えば、図5の上部に示すように、一実施形態において、マネージャー165は、オブジェクトインデックス値を介してすぐにアクセスできるように体系化されるテーブル、又は同様のデータ構造を保存するように設定されてもよい。図示した実施形態において、インデックステーブル500は、多くのエントリ510を含んでもよく、エントリのそれぞれは、オブジェクトインデックスフィールド、オフセットフィールド、オブジェクトサイズフィールド、メタデータサイズフィールド、及び周期的冗長検査(CRC)フィールドを含む多くのフィールドを含んでもよい。図5の下部に示すように、いくつかの例示的なオブジェクト30について、エントリ510のオフセットフィールドは、論理ストレージ空間内の対応するオブジェクト30の開始位置を特定してもよく、オブジェクトサイズ及びメタデータサイズフィールドは、オフセット位置からオブジェクトデータ及びメタデータを拡張する程度を特定してもよい。図示した実施形態において、他の実施形態においてこの順序は反対になる場合があるが、オブジェクトデータは、オブジェクトメタデータの上位である。CRCフィールドは、周期的冗長検査アルゴリズム又は他の適切な種類のチェックサム、又はハッシュアルゴリズムの結果、保存されてもよい。CRCフィールドに最初に保存される値は、オブジェクト30が最初にノード160に保存された時に計算されてもよい。それに続いて、オブジェクト30がアクセスされると、同一のアルゴリズムがオブジェクトデータ及び又はメタデータ、及び保存されたCRCフィールド値と比較した結果値に適用されてもよい。比較が不一致の結果になった場合、保存されたデータの整合性に障害が起きた可能性がある。他の実施形態において、エントリ510は、付加的な、又はこれら示されたものから異なるフィールドを含んでもよいことに留意されたい。例えば、CRCフィールドは省略されるか、又はどこかに実装されてもよい。さらに、オブジェクトデータ及びメタデータの固定位置は、関連のオフセットに加えて、又はそれに代わって保存されてもよい。

20

30

#### 【0072】

リパッカー163は、オブジェクト30が解放される時、及びその関連するストレージリソースが再要求された時に現れるギャップを除くための論理オブジェクトストレージ空間で操作するように設定されてもよい。一実施形態において、リパッカー163は、前回の解放操作によって解放されているSNMコントローラ161及び/又はマネージャー165によってマークされているオブジェクト30を識別するために、論理オブジェクトストレージ空間(例えば、一定期間又は継続的に)スキャンするように設定されてもよい。リパッカー163は、解放されたオブジェクト30の除去を反映するために更新すべき解放されたオブジェクト30のインデックスの後に現れるインデックスを伴うこれらオブジェクト30のエントリ510をもたらしてもよく、これらのオブジェクト30が、論理オブジェクトストレージ空間の元の位置にむかって移動する効果的な結果となる場合がある。例えば、図5の下部のオブジェクトNが解放された場合、リパッカー163は、オブジェクトN+1の新規のオフセットフィールドとして、オブジェクトNのオフセットフィー

40

50

ルドを反映するために更新されるべきオブジェクトN+1と対応するエントリ510をもたらしために操作してもよい。リパッカー163は、削除されるべきオブジェクトNに関連するエントリ510をもたらしてもよく、移動を反映するためにオブジェクトN+1に続くオブジェクトのオフセットを更新してもよい。一実施形態において、マネージャ165は、オブジェクトデータ及びメタデータの対応する移動を、論理オブジェクトストレージ空間及び/又はストレージデバイス169の根本を成すフィールド又は構造内におこさせるように動作してもよい。

#### 【0073】

一部の実施形態において、マネージャ165は、異なる種類のハードウェア及びソフトウェアを含む複数の異なる実行プラットフォーム上で実行するように設定されてもよい。このような一部の実施形態において、抽象概念の1つ以上の付加的な層は、SNMコントローラ161及びそのクライアントにマネージャ165によってもたらされる論理オブジェクトストレージ空間の間に存在してもよい。例えば、図示した実施形態において、マネージャ165は、ファイルシステム167によって管理される1つ以上の物理的ファイルとしての論理オブジェクトストレージ空間を実装するように設定されてもよい。一般的に言うと、ファイルシステム167は、様々な種類の物理的ストレージデバイス169を、本願において物理的ファイルとして言及される論理単位にデータを保存する論理ストレージデバイスに体系化するように設定されてもよい。ファイルシステム167によって管理される論理ストレージデバイスは、事実上、階層的であってもよい。例えば、ファイルシステム167は、物理的ファイルを保存し、それにアクセスするために誘導される場合があるディレクトリ、又はフォルダの階層をサポートしてもよい。一般的にいえば、ファイルシステム167は、所定の物理的ファイルと、物理的ファイルの対応するデータ及び/又はメタデータが保存されるストレージデバイス169の位置間の関係を追跡し、管理するように設定されてもよい。したがって、一実施形態において、マネージャ165は、ファイルシステム167によって配置される1つ以上の物理的ファイルに対する論理的オブジェクトストレージ空間のマッピングを管理してもよい。入れ替わりに、ファイルシステム167は、ストレージデバイス169のアドレス可能な位置に対するこれらの物理的ファイルのマッピングを管理してもよい。

#### 【0074】

根本を成すデバイス169の管理のための異なる特徴を提供する、あらゆる所定のオペレーティングシステムは、様々な異なるファイルシステム167をサポートしてもよいが、ファイルシステム167は、オペレーティングシステム内で概して統合されてもよい。例えば、様々なバージョンのMicrosoft Windows（登録商標）オペレーティングシステムは、FAT32（File Allocation Table - 32）及びFAT16ファイルシステムのみならず、NTファイルシステム（NTFS）などのファイルシステムをサポートする。様々なバージョンのLinux and Unix（登録商標）オペレーティングシステムは、ext/ext2ファイルシステム、Network File System（NFS）、Reiser File System（ReiserFS）、Fast File System（FFS）、その他多数のファイルシステムをサポートしてもよい。一部の第三者ソフトウェアベンダーは、例えば、VERITAS（登録商標）File System（VxFS）などの様々なコンピューティングプラットフォームとの統合のための独自のファイルシステムを提供してもよい。異なるファイルシステムは、根本を成すストレージデバイス169の管理のため、様々な特徴のためのサポートを提供してもよい。例えば、一部のファイルシステム167は、デバイスミラーリング、ストライピング、スナップショットティング、又は他の種類の仮想化特徴のためのサポートを提供してもよい。

#### 【0075】

一部の実施形態において、マネージャ165とストレージデバイス169間に依然として抽象概念のさらなる層が存在してもよい。例えば、一部の実施形態において、ボリュームマネージャ層は、ファイルシステム167とストレージデバイス169の間に提供

10

20

30

40

50

されてもよく、上述の仮想化特長のいくつか、あるいはすべてを行うように設定されてもよい。代わりに、特定のストレージデバイス169は、ハードディスクドライブの独立した配列、又は仮想化コントローラを含む他のデバイスとして設定されてもよい。上記に記載されるように、ボリュームマネージャーによって、又はファイルシステム167内でサポートされてもよい仮想化マッピングと同様に、内部で、仮想化コントローラは、ディスクドライブヘドバースのストレージアドレス空間の複雑なマッピングを任意にサポートしてもよいが、仮想化コントローラは、単一の物理的デバイスとして、ファイルシステム167にディスクドライブをもたらすように構成されてもよい。一部の実施形態においては、示した抽象概念の層よりも少ない場合があることに留意されたい。例えば、一部の実施形態において、マネージャー165は、ファイルシステム167を使用することなく、例えば、未加工の物理的デバイスとしてのストレージデバイス169と直接交流するように設定されてもよい。

#### 【0076】

一般的に言えば、ストレージデバイス169は、ファイルシステム167及び/又はマネージャー165によってサポートされてもよいストレージデバイスのあらゆる適切な種類を含んでもよい。ストレージデバイス169は、Small Computer System Interface (SCSI) デバイス、又はAT Attachment Programming Interface (ATAPI) デバイス (Integrated Drive Electronics (IDE) デバイスとしても知られる) のようなハードディスクドライブデバイスを一般に含んでもよい。しかしながら、ストレージデバイス169は、磁気、又は光学媒体系のデバイス、ソリッドステート大容量記憶装置 (例えば、不揮発性の、又は「フラッシュ」メモリ系のデバイス)、磁気テープなどを含む、あらゆる種類の大容量記憶デバイスを網羅してもよい。さらに、ストレージデバイス169は、Universal Serial Bus又はIEEE 1394/Firewire (登録商標) 基準のバージョンに準拠するインターフェースなど、上述のものに加え、あらゆる適切なインターフェース種を通じてサポートされてもよい。

#### 【0077】

上記に記載されるように、ストレージサービスシステム内に保存されるオブジェクト30のあらゆる所定のインスタンスについて、対応するロケータは、該システム内のすべてのノード160にわたるインスタンスを個別に識別してもよい。一実施形態において、ロケータは、個別の識別子、又はオブジェクトインスタンスが保存されるノード160と対応する、「ノードID」のみならず、マネージャー165によってオブジェクトに指定される、連結、組み合わせ、又は、オブジェクトインデックス値の他の機能として作成されてもよい。例えば、上記に記載されるように、64ビットオブジェクトインデックス値は、128ビットロケータを生じる64ビットノードIDと組み合わせてもよい。様々な実施形態において、ロケータを形成するために、より少ない、又は大きい数のビットが用いられてもよいが、このようなロケータは、 $2^{64}$ と同数の固有ノード160が $2^{64}$ と同数の固有のオブジェクトインスタンスを保存することを可能にする。

#### 【0078】

一実施形態において、ノードIDは、タイムスタンプ、又はデートスタンプを有する、所定のノード160と対応する、Internet Protocol (IP) アドレスなどの、固有のネットワークアドレスの連結、又は結合を通じて形成されてもよい。例えば、ノード160は、IPアドレスが指定された時点とを反映するタイムスタンプと、又はIPアドレスが有効であるとわかっている間の時点との結合において、そのIPアドレスに従い (例えば、ノードスタートアップ/初期設定において、又はノードIDが指定された時点で、さもなければ、初期設定の間に)、ノードIDを指定されてもよい。一般的に言えば、同一のIPアドレス空間に所属する2つの相異なるノード160は、いかなる所定の時において、同一のIPアドレスを有効に指定されることはない。したがって、ノードのIPアドレスとタイムスタンプの値の結合は、そのノードに固有の識別子を生じてもよい。例えば、32ビットのIPアドレスは、他のビット幅も用いられてもよいが、上述の

10

20

30

40

50

64ビットノードIDを生じるために、32ビットのタイムスタンプ（例えば、一部の共通の参照時間のため、それは経過した秒数を示す）と連結、又は結合されてもよい。ノードIPアドレスに頼らない固有のノードIDを指定するために、他の技術が使用されてもよいことを考慮されたい。例えば、ネームサーバなどの中央局は、上記に記載されるように、ノード160内のオブジェクトインデックス値の指定と同様に、ノードIDの固有性を保証する方法での要求におけるノードIDを委任してもよい。

#### 【0079】

ノードIDがノードIPアドレスから派生したノードIDで一部の実施形態において、該ノードIDは、いかなる時にもノード160の現在のIPアドレスを反映しなくてもよいことに留意されたい。例えば、ノードIDは、ノード160が初期化されるまで存続するが、ノードのIPアドレスは、ノードIDの作成に続いて変更されるか、又は再指定されてもよい。また、一部の実施形態において、ノードIDは、ストレージクライアント50又は他の悪意のあるエンティティが実際のノードIPアドレスを決定するためにロケータを解読することから防ぐために、決定的な方法でハッシュされ、暗号化され、又は分かりにくくなる場合がある。

#### 【0080】

図4のノード160の実施形態に関する取得、格納、及び解放操作の例示的な実施形態の操作は、図6から図8に示される。最初から図6までに言及して、取得操作は、該操作がコーディネータ120又は他のクライアントからノード160で受信される、ブロック600で始まってよい。例えば、コーディネータ120は、上記に記載されるように、ノードID及びオブジェクトインデックス値を含む特定のロケータに取得操作を発行してもよい。該ノードIDは、例えば、ノードIDが目標とするノード160の現在のIPアドレスを反映する場合、適切なノード160に対して取得操作を直接送るために使用されてもよい。代わりに、以下に説明するように、DFDD110などのディレクトリサービスが、ノードIDを、取得操作が適切なノード160に送られることを通じるアドレス可能なエンドポイント又は目的地に分解するために用いられてもよい。

#### 【0081】

ノード160によって受信されると、取得操作は、ノード160の論理オブジェクトストレージ空間内の目標とされるオブジェクトインスタンスの範囲を識別するために処理されてもよい（ブロック602）。例えば、コントローラ161は、取得操作を受信し、それをマネージャー165へ伝えてもよい。入れ替わりに、マネージャー165は、論理オブジェクトストレージ空間内の所望のオブジェクトインスタンスの位置を得るために取得操作がインデックステーブル500にアクセスすることによって参照される、ロケータのオブジェクトインデックス部分を使用してもよい。例えば、マネージャー165は、オフセットからオブジェクトインスタンスの長さのみならず、オブジェクトインスタンスが開始する、論理オブジェクトストレージ空間にオフセットを得てもよい。一部の実施形態において、取得操作は、オブジェクトデータ、メタデータ、又はその両方が要望されるかを特定してもよい。そのような実施形態において、マネージャー165は、要求データと関連し、論理オブジェクトストレージ範囲を決定してもよい。例えば、オブジェクトデータ及びメタデータの両方が要望される場合、マネージャー165は、取り出すべきオブジェクトオフセットからの範囲を決定するために、オブジェクトデータサイズ及びメタデータサイズの両方を使用してもよい。上記に述べたように、他の実施形態において、オブジェクトインスタンスのためのストレージ範囲は、論理オブジェクトストレージ空間内の関連するオフセットよりは、固定位置を通じるなどして、異なる方法でマネージャー165によって保存、及び管理されてもよい。

#### 【0082】

論理オブジェクトストレージ空間内のオブジェクト範囲は、物理的ファイルストレージ空間内の1つ以上の対応するファイル内のエクステンツにマップ化されてもよい（ブロック604）。例えば、マネージャー165は、論理オブジェクトストレージ空間を、ファイルシステム167によって管理される1つ以上のファイルにマップし、例えば、読み込

10

20

30

40

50

まれるべき指定されたファイル内のロケータ、又はオフセットのみならず、1つ以上のファイルネームを反映することによって、所望のオブジェクト範囲と対応するデータを得るために、ファイルシステム167に適切なファイルアクセス操作を発行してもよい。他の実施形態において、コントローラ161は、マネージャ165によって管理される論理ブロックストレージ空間を迂回するように設定されてもよく、代わりに、ファイルシステム167によって管理される物理ファイルと直接交流してもよいことを考慮されたい。

#### 【0083】

物理的ファイルへの参照は、デバイス関連要求(ブロック606)に対してマップ化されてもよい。例えば、ファイルシステム167は、論理ブロックアドレス(LBA)又はデバイス形状(例えば、シリンダ、トラック、セクタ及び/又はヘッド)に特定のアドレスなど、ストレージデバイス169のアドレス可能な位置を特定するための1つ以上の読み込み要求を作成するように設定されてもよい。上記に述べたように、一部の実施形態において、マネージャ165は、ファイルシステム167を迂回し、ストレージデバイス169を直接管理するように設定されてもよい。

#### 【0084】

要求されたオブジェクトデータは、ストレージデバイス169(ブロック608)から取り出され、要求するクライアント(ブロック610)へ返されてもよい。例えば、取り出されたデータは、図4に階層的に示される要求を通じたバックアップを通過されてもよく、又は、要求するクライアントへの輸送のために、ストレージデバイス169、又はファイルシステム167からコントローラ161へ直接返されてもよい。

#### 【0085】

図7に示すように、一実施形態において、図6のブロック600について上記の記載と同様の方法で、操作が、コーディネータ120、又は他のクライアントからノード160で受信される時、格納操作はブロック700で開始されてもよい。例えば、コーディネータ120は、ノードピッカー130によって作成される書き込みプランにおいて特定されるノード160に対して格納操作を発行してもよい。取得操作とは対照的に、格納操作は、保存されるべきオブジェクトデータ及び/又はメタデータを含んでもよく、任意に、データ及び/又はメタデータの長さを指定する付加的なパラメータを含んでもよい。

#### 【0086】

ノード160によって受信されると、格納操作は、論理オブジェクトストレージ空間内のオブジェクトインスタンスのためのストレージ範囲を指定するために処理されてもよい(ブロック702)。一実施形態において、マネージャ165は、オブジェクトインデックス値を新規オブジェクトインスタンスに指定し、インデックステーブル500を新規オブジェクトインスタンスのオフセットを特定する新規エントリ510に記録するように設定されてもよい。例えば、新規エントリのオフセットは、インデックスの最高値を有する存在するオブジェクトインスタンスのストレージ範囲(例えば、オフセット及び長さ)に関連して決定されてもよい。新規オブジェクトインスタンスのデータ及び/又はメタデータの長さが格納操作に対するパラメータとして特定されない場合、マネージャ165、又はコントローラ161は、新規エントリ510における含有としてこれらを計算するように設定されてもよい。

#### 【0087】

論理オブジェクトストレージ空間内の新規に指定されたストレージ範囲は、物理的ファイルストレージ空間内の1つ以上の対応するファイル内のエクステンツに対しマップ化されてもよい(ブロック704)。例えば、新規オブジェクトインスタンスに対して指定された範囲は、1つ以上の存在する物理的ファイルの最後に加えられてもよく、そうでない場合は、存在する、又は新規に割り振られた物理的ファイル内に位置されてもよい。物理ファイル範囲は、例えば、取得操作について上記の記載と同様の方法で、ファイルシステム167によってストレージデバイス範囲にマップ化されてもよく(ブロック706)、オブジェクトインスタンスデータ及び/又はメタデータは、ストレージデバイス169に保存されてもよい(ブロック708)。

## 【 0 0 8 8 】

データ及びノ又はメタデータがうまくストレージデバイス 169 に書き込まれたことを確認すると、保存されたオブジェクトインスタンスに対応するロケータが、要求するクライアントに返されてもよい（ブロック 710）。例えば、マネージャー 165 は、ノード 160 のノード ID に作成されたオブジェクトインデックス値を加えるように設定されてもよく、物理的ファイル書き込み操作が無事に完了するファイルシステム 167 からの表示に際しオブジェクトロケータとして結果値を返してもよい。

## 【 0 0 8 9 】

図 8 に示すように、一実施形態において、図 6 のブロック 600 について上記の記載と同様の方法で、操作がコーディネータ 120、又は他のクライアントからノード 160 で受信されるとき、解放操作は、ブロック 800 で開始されてもよい。解放された操作は、他の実施形態において、他の引数も提供されるが、解放されるべきオブジェクトインスタンスのロケータを単に特定してもよい。

## 【 0 0 9 0 】

取得操作のように、ノード 160 によって受信されると、解放操作は、ノード 160 の論理オブジェクトストレージ空間内の目標とされるオブジェクトインスタンスの範囲を識別するために処理されてもよい（ブロック 802）。例えば、コントローラ 161 は、解放操作を受信し、それをマネージャー 165 へ伝えてもよい。入れ替わりに、マネージャー 165 は、参照されるオブジェクトインスタンスの対応するエントリ 510 を識別するために、解放操作がインデックステーブル 500 にアクセスすることによって参照される、ロケータのオブジェクトインデックス部分を使用してもよい。参照されるオブジェクトは、解放されたとしてマークされてもよい（ブロック 804）。例えば、マネージャー 165 は、エントリがもはや有効ではないことを表す、負の数のような不正値に対し、オフセット、又は別のフィールドエントリ 510 を設定するように設定されてもよい。肯定応答が、オブジェクトが解放されたことを示す要求するクライアントに返されてもよい（ブロック 806）。

## 【 0 0 9 1 】

上記に記載されるように、オブジェクトインスタンスと関連するストレージリソースは、オブジェクトインスタンスが解放された時に他の使用のために直ちに解放、再要求、又は再配分されてもよい。むしろ、一実施形態において、これらの供給源は、解放操作がこれらを再要求することに関し、非同期的に操作する独立過程まで存続してもよい。図 9 は、例えばストレージリパッカー 163 によって実装されてもよいような、当該の過程の一実施形態の操作を示す。ブロック 900 において、ノード 160 に保存された、特定のオブジェクトインスタンスに対応するオブジェクトインデックスエントリが選択されてもよい。例えば、リパッカー 163 は、エントリに保存されたオブジェクトインデックス値に従い、配列の順序で、インデックステーブル 500 からインデックスエントリ 510 を選択するように構成されてもよい。それに続いて、選択されたエントリは、対応するオブジェクトインスタンスが解放されたかどうかを決定するために調査されてもよい（ブロック 902）。例えば、リパッカー 163 は、フィールドが、負の数値、又はある他の値などの、対応するオブジェクトインスタンスが解放されたことを示す値を設定したかどうかを確かめるために、オフセットフィールド、又は別のフィールドを調べてもよい。

## 【 0 0 9 2 】

選択されたオブジェクトが解放されていない場合、操作は、別のオブジェクトが選択されるブロック 900 に戻って続行されてもよい。選択されたオブジェクトが解放された場合、論理オブジェクトストレージ空間は、解放されたオブジェクト（ブロック 904）と対応するストレージリソースに再要求するために、再圧縮されてもよい。例えば、リパッカー 163 は、第 1 の当該オブジェクトインスタンスのオフセットが解放されたオブジェクトのオフセットに設定され、次の当該オブジェクトインスタンスのオフセットがデータサイズの機能として設定され、メタデータサイズと第 1 の当該オブジェクトインスタンスのオフセット、など、論理オブジェクトストレージ空間内の解放されたオブジェクトに続

く、これらのオブジェクトインスタンスのインデックスエントリ 510 を調整するように設定されてもよい。しかしながら、一部の実施形態において、解放されたオブジェクトインスタンスに続くオブジェクトインスタンスのすべてが、新規オブジェクトが調査のために選択される前に再圧縮される必要があるわけではない。例えば、再圧縮は、接触する各オブジェクトが、それが調査のために選択される時に圧縮されるように、オブジェクト選択と交互配置してもよい。

#### 【0093】

一部の実施形態において、マネージャー 165 は、論理オブジェクトストレージ空間の再圧縮に応じて物理的ファイルストレージ空間内で、同様の再圧縮、又は圧密操作を行ってもよい。例えば、マネージャー 165 は、論理オブジェクトデータ範囲が、異なる物理的ファイルデータ範囲に再マップ化されることを生じてもよい。同様に、一部の実施形態において、ファイルシステム 167 は、物理的ファイルストレージ空間の再圧縮に応じてストレージデバイス 169 間で、類似の再圧縮、又は圧密操作を行ってもよい。他の実施形態において、物理的ファイルストレージ空間、又はストレージデバイスそれ自体の再圧縮は、リバッカー 163 によって開始される論理オブジェクトストレージ空間再圧縮を独立して発生してもよい。例えば、ファイルシステム 167 は、マップ化されたデバイスストレージエクステンツが、ストレージデバイスのアクセス傾向に関連して、ほとんど、又は完全に連続的であるような、デバイスストレージエクステンツへの物理的ファイルストレージエクステンツのマッピングを再配置することによって、ストレージデバイスに保存された物理的ファイルを最適化するように構成されてもよい。

#### 【0094】

論理オブジェクトストレージ空間の圧縮に続き、解放されたオブジェクトに対応するインデックスエントリが削除され（ブロック 906）、操作は、別のオブジェクトが選択されるブロック 900 から継続されてもよい。上記に述べたように、一部の実施形態において、複数オブジェクトが選択された場合「オンザフライ」で再圧縮が発生してもよく、オブジェクトを再配置するために要求される操作の数を最小限にするため、論理オブジェクトストレージ空間の全体的な稼働率を改善してもよい。

#### 【0095】

一部の実施形態において、ノード 160 によってサポートされてもよい取得、格納、解放、又は他の操作のいずれかは、要求するクライアントに関する様々な種類のハンドシェイク、肯定応答、又はエラー処理プロトコルをサポートしてもよいことに留意されたい。例えば、クライアントが操作に対し、不正な形式での要求を要求する場合（例えば、必要なパラメータを供給することを怠る）、又はノード 160 が、操作を十分に完了できない場合（例えば、それが格納操作を引き受けるのに不十分な供給源を有する）、ノード 160 が要求するクライアントにエラー表示を返してもよい。このような表示は障害状態の性質のため、特定の詳細を含む、又は含まなくてもよい。

#### 【0096】

一実施形態において、コーディネータ 120 は、複数の操作が共通のデータを有してもよいが、操作によって目的とされる各それぞれのノード 160 に対する操作を独立して伝えるように設定されてもよい。例えば、書き込みプランに従い、複数のノード 160 にオブジェクト 30 が書き込まれる格納操作の事例において、コーディネータ 120 は、各特定のノード 160 と独立して通信してもよい。しかしながら、他の実施形態において、複数目的地ノード 160 に意図された共通データ及び/又はパラメータを有する操作は、連鎖されてもよい。一実施形態において、コーディネータ 120、又は他のクライアントは、受信者リストなど、操作のパラメータにおける各受信者を特定することによって、連鎖操作を先導してもよい。操作において示唆される複数の受信者は、初期設定によって連鎖を示してもよく、又は別のパラメータが、操作が連鎖されていることをマークするために使用されてもよい。コーディネータ 120、又は他のクライアントが、操作に特定される最初の目的地ノード 160 の 1 つに連鎖操作を伝えることによって先導してもよい。

#### 【0097】



連鎖操作を受け取る際に、ノード160は、操作を処理し、それを操作において特定される別の目的地ノード160の1つに転送してもよい。当該の転送に先立ち、受け手ノード160は、受け手を示し、循環転送を避けるために、操作に含まれる目的地リストからそれ自体を取り除いてもよい。操作は、受け手ノードの処理と同時に転送されてもよい。代わりに、転送は、受け手ノードの無事な処理の完了を条件としてもよい。一部の実施形態において、連鎖操作は、受信者が操作内で示される順序で受信者に伝えられてもよい。他の実施形態において、ノード160は、例えば、残っている目的地のどれが一番近いのか、最低負荷であるか、又は、いくつかの他の選択基準を満たすかを決定することによって、動的に次の受信者を選択してもよい。一部の実施形態において、連鎖、及び非連鎖操作の組み合わせは、コーディネータ120、又は他のクライアントによって作成されてもよいことに留意されたい。例えば、同一のデータが、6つの相異なるノード160行きである格納操作の目標である場合、コーディネータ120は、6つの目的地ノード、又はそれぞれが3つの目的地ノードを特定する2つの連鎖操作を特定する、単一連鎖操作を作成してもよい。コーディネータ120が、各それぞれの目的地ノード160に独立して伝える、6つの非連鎖操作の作成を含む、他の組み合わせも可能である。

10

#### 【0098】

##### キーマップ構造

上記に記載されるように、様々なビットストアノード160は、オブジェクト30のインスタンスのためのストレージを提供するように設定されてよい。ノード160は、個別に冗長、又はデータセキュリティのためのあらゆる特定のサポートを提供しなくてもよく、実際に、一部の実施形態において、ノード160は、オープンソースのオペレーティングシステム（例えば、Linux）を作動する包括的なコンピュータプラットフォームを使用し、安価な、商品ハードドライブ（例えば、ATAPI/DEハードドライブ）を介してストレージを提供することで実装されてもよい。このような実施形態において、個々のシステムは特にフォルトトレラントでなくてもよい。むしろ、データセキュリティ及び冗長は、上記に記載されるように、多くのノード160に渡るオブジェクト30の複製を通じて提供されてもよい。

20

#### 【0099】

前に論述したように、所定のオブジェクト30は、ストレージクライアントによって特定されてもよいキーと対応してもよい。所定のオブジェクト30の個々のインスタンスは、ストレージサービスシステム内に含まれるノード160の収集物にわたり、これらインスタンスを固有に識別してもよいそれぞれのロケータに対応してもよい。一実施形態において、ストレージサービスシステム内に展開される各キーマップインスタンス140は、所定のオブジェクト30のためのキー及びすべての対応するロケータと、ノード160間に保存されたその複製されたインスタンスとの間の関係、又はマップを保存し、維持するように設定されてもよい。以下の論考において、どのようにキーマップインスタンス140の特定の実施形態が実装されてもよいかの説明に続き、キーマップインスタンス140の様々な実施形態の一般的な特徴及び機能性を論述する。

30

#### 【0100】

一実施形態において、所定のキーマップインスタンス140は、1つ以上のテーブル、又はあらゆる他の適切な種類のデータ構造内の様々なキーと関連するロケータ間の関係の詳細を保存するように設定されてもよい。例えば、図10に示す一実施形態において、キーマップインスタンス140は、多くのエントリ144を有するキーマップデータ構造142を含む。各エントリは、関連する記録148のみならず、それぞれのキー146を含む。一部の実施形態において、以下にさらに詳しく説明するように、エントリ144を体系化するために使用されるデータ構造の体系は複雑であってもよい。しかしながら、機能的な見地から、キーマップインスタンス140は一般的に、所定のキー144とその対応する記録148との間に、1対1のテーブル型の関係を維持してもよい。

40

#### 【0101】

記録148は、概して、所定のキー144と対応するロケータを含んでよいが、他の情

50

報を含んでもよい。例えば、記録 1 4 8 の一実施形態は以下のように設定されてもよい。

```
struct KeyRecord {
    int16_t version;
    int16_t storageClass;
    int64_t creationDate;
    int64_t objectSize;
    uint32_t crc32;
    int8_t numLocators;
    struct locator {
        int64_t nodeID;
        int64_t objectIndex;
    } replicas [];
}
```

10

このデータ構造例は、C プログラミング言語の構文を使用して表現されるが、あらゆる適切な言語、表現、又は形式を使用して実装されてもよい。記録 1 4 8 の他の実施形態は、これら示されたものよりもさらに多い、さらに少ない、又は異なるフィールドを含んでもよい。一部の例では、記録 1 4 8 は、一部の種の U n i x ファイルシステムに用いられるアイノード構造に対するストレージ空間の体系化における記録 1 4 8 の相似の目的で示す、「アイノード」といわれることもある。しかしながら、該文脈中の「アイノード」という用語の使用は、ファイルシステム、又は他のストレージ文脈内のアイノードの実装、又は使用の特定の詳細を引き合いに出すことを意図しない。

20

#### 【 0 1 0 2 】

上記の実施形態において、記録 1 4 8 は、7 つの特定の要素を含む。1 6 ビットバージョンの要素は、記録 1 4 8 の形式に特定である、個別の識別値を保存するために使用されてもよい。例えば、異なる種類の記録 1 4 8 は、キーマップインスタンス 1 4 0 の異なる実装に使用されてもよく、一部の実施形態において、所定のキーマップインスタンス 1 4 0 内に保存される記録 1 4 8 は、異種であってもよい。バージョン要素は、記録 1 4 8 の異なるバージョン間で区別するために使用されてもよく、記録の他の要素が適切に解読され、使用されてもよい。

#### 【 0 1 0 3 】

30

1 6 ビット `storageClass` 要素は、記録 1 4 8 と対応するオブジェクト 3 0 のストレージクラスの表示を保存するために使用されてもよい。ストレージクラスは、次の項にさらに詳しく述べる。一般的に言えば、オブジェクトの所定のストレージクラスは、所定のストレージクラスの他のメンバーに共通であってもよいストレージ特徴及び/又はポリシーを識別してもよいが、他のストレージクラスのメンバーからは異なってもよい。例えば、「高信頼性」ストレージクラス及び「低信頼性」ストレージクラスは、ストレージサービスシステムの所定の実装のために定義されてもよい。高信頼性ストレージクラスのメンバーであるオブジェクト 3 0 は、低信頼性ストレージクラスのメンバーであるオブジェクト 3 0 のより高い程度にたいして複製されてもよく、したがって、おそらく、低信頼性ストレージクラスのメンバーにたいして見積もられたものよりも高い使用コストと引き換えに、個々の複製の損失に対する敏感性を減少してもよい。多くの他の可能なストレージクラスの種類及び組み合わせが可能であり、考慮される。

40

#### 【 0 1 0 4 】

6 4 ビット `creationDate` 要素は、対応するオブジェクト 3 0 がストレージサービスシステム内に作成された日付及び時間の表示を保存するために使用されてもよい。この要素は、あらゆる適切な方法に形式化されてもよい。例えば、日付及び時間は、共通の基準点のため、要素内の相異なるフィールドとして、又は経過時間単位の数（例えば、秒、ミリ秒など）を表す単一の数として明示的にコード化されてもよい。一部の実施形態において、最近の修正要素は、他の実施形態において、記録 1 4 8 内の相異なる要素として含まれてもよいが、`creationDate` 要素は、対応するオブジェクト 3 0 の

50

いかなる局面の最近の修正の日付及び時間を示すように構成される、付加的なフィールドを含んでもよい。

【 0 1 0 5 】

64ビットobjectSize要素は、例えば、バイトで、対応するオブジェクトのサイズの表示を保存するために使用されてもよい。一部の実施形態において、この要素は、他の実施形態において、これらは相異なるフィールドとして保存されてもよいが、オブジェクトデータ及びメタデータの両方のサイズを反映してもよい。32ビットcrc32要素は、あらゆる適切なチェックサムアルゴリズムに従い、オブジェクトデータ及び/又はメタデータのための計算された周期的冗長検査(CRC)チェックサムの表示を保存するために使用されてもよい。例えば、チェックサムは、変造、又は改ざんに対するデータ整合性を立証するために含まれてもよい。他の実施形態において、オブジェクトデータ及び/又はメタデータから計算されたあらゆる適切な種類の、ハッシュ、又は署名は、CRCチェックサムに加えて、又はそれに代わって使用されてもよい。

10

【 0 1 0 6 】

8ビットnumLocators要素は、複製[ ]配列内の記録148内に含まれるロケータの数の表示を保存するために使用されてもよい。この配列内で、各ロケータは、64ビットオブジェクトインデックス値のみならず、64ビットノードID要素として保存され、ビットストアノード160の構造の論考で上記に記載されるように作成されてもよい。一部の実施形態において、ロケータは、複製[ ]配置内に単一要素として保存されてもよい。

20

【 0 1 0 7 】

一実施形態において、キーマップインスタンス140は、提供されるAPIをサポートするために必要なこれらの機能を実行するのみならず、コーディネータ120などの、キーマップクライアントに対してキーマップAPIを提供するように設定されてもよい。例えば、コントローラ120は、キーマップインスタンス140によって管理されるエントリ144と関連する記録148上で、保存、取り出し、削除、又は他の操作を行うために、APIを使用するように設定されてもよい。上記に記載されるようにノード160によってサポートされてもよいオブジェクトインスタンス上の操作と類似して、一実施形態において、キーマップAPIは、キーマップエントリ144上の格納、取得、及び削除操作をサポートしてもよい。このような一実施形態において、キーマップ保存操作、又はキーマップ書き込み操作として総称的に言及されてもよい、キーマップエントリ格納操作は、キーマップエントリ144内に保存されるべきキー146及び記録148を特定してもよい。一実施形態において、エントリ144が既に存在するキー146を特定する格納操作は、存在するエントリ144に関連する記録148を、格納操作の引数、又はパラメータとして特定される記録と置き換えてもよい。所定のキーマップインスタンス140が完了次第、キーマップ格納操作は、例えば、操作が成功、又は失敗したか、任意の種類の障害が生じたか(もしある場合)などの、状態表示を要求するクライアントに返してもよい。一部の実施形態において、キーマップ格納操作が存在するエントリ144の再配置の結果になった場合、キーマップインスタンス140は、要求するクライアントに対するエントリ144の前回値を返すように設定されてもよい。

30

40

【 0 1 0 8 】

総称的にキーマップ読み込み、又は取り出し操作といわれることもあるキーマップエントリ取得操作は、一実施形態において、パラメータとしてキーを特定してもよい。完了次第、キーマップ取得操作は、当該のエントリが存在する場合、要求されるキーに関連するキーマップエントリ144の記録148を要求するクライアントに返してもよい。対応するエントリ144が存在しない場合は、その要旨が要求するクライアントに返されてもよい。

【 0 1 0 9 】

一実施形態において、キーマップエントリ削除操作は、要求するクライアントが、エントリに記録を書き込むことを特定する必要がない場合を除き、格納操作と同様に行われる

50

ように設定されてもよい。所定のキーマップインスタンス 140 が完了次第、キーマップ削除操作は、キーマップ格納操作と同様の状態表示を要求するクライアントに返してもよい。格納操作のように、一部の実施形態において、キーマップインスタンス 140 は、要求するクライアントに対する削除されたエントリ 144 の前回値を返すように設定されてもよい。

#### 【0110】

キーマップAPIは、様々な実施形態において、他の種類の操作をサポートしてもよい。例えば、キーマップAPIは、キーマップエントリの管理において、キーマップクライアントを援助する操作をサポートしてもよい。一実施形態において、キーマップAPIは、要求するクライアントによって特定されるいくつかの基準と一致する、キー 146 を有するこれらのエントリ 144 を識別するように設定されてもよいリスト操作を支持してもよい。例えば、リスト操作は、クライアントが、操作に対するパラメータとしての文字列、又はパターンを特定できてよい。所定のキーマップインスタンス 140 が完了次第、リスト操作は、要求するクライアントに、特定された文字列、又はパターンを満たすこれらのキー 146 のリストを返してもよい。一実施形態において、キー 146 は、キー 146 の適切なプレフィックスである場合にのみ、所定の文字列を満たしてもよい（例えば、文字列のすべての文字について、文字列のN番目文字がキーのN番目文字と一致する）。他の実施形態において、キー 146 は、文字列が、キー 146 内のいかなる場所においても見つかる場合、所定の文字列を満たしてもよい。

#### 【0111】

リスト操作は、一部の実施形態において、他のパラメータをサポートしてもよい。例えば、リスト操作は、要求するクライアントが、返されるべき一致数の限度を特定出来てもよい。さらに、要求するクライアントは、例えば、検索すべきキー 146 が当てはまるオープンエンド、又はクローズエンドの辞書式範囲内で、検索するキー 146 に対する制約を特定してもよい。一部の実施形態において、キーマップインスタンス 140 は、リスト操作基準を満たすキー 146 のみならず、記録 148 を返すように設定されてもよい。また、一部の実施形態において、キーマップAPIは、リスト操作として、同一の種類のパラメータ及び実行動作をサポートしてもよい、カウント操作をサポートしてもよい。しかしながら、要求するクライアントによって提供される基準を満たすこれらのキー 146 及び/又は記録 148 を返す代わりに、カウント操作は、これらの基準を満たすキーの数を返してもよい（例えば、対応するリスト操作によって返されるであろうキーの数）。またキーマップAPIは、上記に説明されない他の操作もサポートしてもよいことに留意されたい。

#### 【0112】

一部の状況において、異なるキーマップクライアントは、同一のキーマップエントリ 144 を修正することを求める場合がある。例えば、様々なクライアント、又はシステムによる操作に応じ、2つの異なるコーディネータ 120 は、所定の記録 148 の内容を同時に変更することを試みてよく（例えば、複製のロケータを追加、削除、又は修正するために）、又はもう一方が対応するエントリ 144 を削除することを試みる一方、1つが記録 148 を修正することを試みてよい。所定のキーマップエントリ 144 への同時要求を解決するための一貫した方法を提供するために、一実施形態において、キーマップAPIは、キーマップ状態を更新、又は修正する（例えば、キーマップ格納及び削除操作）少なくともこれらのキーマップ操作が、キーマップ操作に対するパラメータとしての配列番号を提供することを必要としてもよい。キーマップインスタンス 140 は、配列番号を比較すること（例えば、数的に、又は辞書編集的に）、及び比較に基づく操作の1つを継続的に選定することによって、エントリ 144 に対する相反する更新を解決するように設定されてもよい。一部の実施形態において、以下にさらに詳しく説明するように、同期化の回復のために、提供される配列番号は、修正された記録 148 と共に、修正されたキーマップエントリ 144 に保存されてもよい。

#### 【0113】

例えば、キーマップクライアントは、タイムスタンプに基づいて配列番号を作成してもよい。一実施形態において、当該のタイムスタンプは、以下のように形式化された64ビット数を含んでもよい。タイムスタンプのビット63は、ゼロに設定されてもよい（例えば、タイムスタンプが署名された、又は署名されない番号であるかどうかの混乱を避けるため）。ビット62:32は、参照時から経過した秒の数を含んでもよい（例えば、1970年1月1日午前零時、グリニッジ標準時、Unix及びLinuxの多くのバージョンによって用いられる参照時間）。ビット31:22は、最終秒から経過したミリ秒の数を含んでもよい。ビット21:0は、実質的に無作為に作成されたビットを包含してもよい。他の実施形態において、タイムスタンプは、フィールドの異なる幅、又は種類に基づいて作成されてもよい。代わりに、キーマップクライアントは、配列番号を作成するための、全く異なる基準を用いてもよい。提供される配列番号の解像度が高ければ、同一のキーマップエントリ144に対する異なるキーマップクライアントによって提供される、異なる配列番号間のコリジョンの可能性は低い場合がある。しかしながら、コリジョンが発生する場合、キーマップインスタンス140は、あらゆる適切な、一貫した技術を使用してコリジョンを解決するように設定されてもよい。

#### 【0114】

多くの実施形態において、マッピングキーにおいて、ロケータに対するキーマップインスタンス140の抽象的機能性の動作は、比較的直接的である。例えば、上記に記載されるように、キーマップインスタンス140の一実施形態によってサポートされる基礎的な操作の一式は、記録148内に含まれるキー146とロケータとの間の関係を反映するエントリ144を操るよう構成される、格納、取得、及び削除操作を含んでもよい。しかしながら、ストレージサービスシステム内のキーマップ機能性の実装は多くの困難をもたらす場合がある。とりわけ、ストレージサービスシステムが、多数のクライアントに代行して、膨大な数のオブジェクト30をサポートする場合（例えば、ストレージの合計でテラバイト（TB）、又はペタバイト（EB）になる数百万、又は数十億、又はそれ以上のオブジェクト30）、同様に容量を拡大するために、キーマップの実装が必要となる場合がある。しかしながら、単独のコンピュータシステム内のキーマップに含まれる全情報を示す、十分なシステムメモリ供給源を実装することは可能でない、又は経済的に実現可能でない場合がある。さらに、フォルトトレランス及び増加するキーマップクライアント要求のための処理スループットについて、キーマップデータの複数の複製が、ストレージサービスシステム内に分散される方法で展開されてもよい。しかしながら、キーマップデータの複製は、例えば、1つの複製が修正される一方、もう1つがアクセスされる場合、キーマップ同期化及び一貫性問題につながる場合がある。

#### 【0115】

キーマップ機能性の拡張性は、キーマップインスタンス140内の階層のレベルを導入することによって改善されてもよい。このような階層の一実施形態は、図11A～Dに示す。図11Aにおいて、例のキーマップ配置1100が示される。上記に記載されるように、例えば、図3に関し、一部のストレージサービスシステムの実施形態において、複数のキーマップインスタンス140が、例えば、異なるデータセンタ300又は領域310に、システム全体にわたり分散されてもよい。概して、キーマップインスタンスの収集物は、展開といわれることもある。一部の実施形態において、他の実施形態においては、システムは、キーマップの階層の付加的なレベル下に統合される複数のキーマップ展開1100を含んでもよいが、ストレージサービスシステムは、システム内に支給されるすべてのキーマップインスタンス140を含む、単一のキーマップ展開1100を網羅してもよい。

#### 【0116】

図示した実施形態において、展開1100は、キーマップインスタンス140a～cを含み、例えば、以下にさらに詳しく述べるようにインスタンス同期プロトコルに従い、それぞれは他とキーマップ情報を交換するように構成される。示すように、各キーマップインスタンス140は、互いに通信しあうように構成される多くのホスト400を含む。例

10

20

30

40

50

えば、キーマップインスタンス 140 a は、ホスト 400 a ~ c を含み、キーマップインスタンス 140 b は、ホスト 400 d ~ g を含み、キーマップインスタンス 140 c は、ホスト 400 h ~ j を含む。一般的に言えば、各ホスト 400 は、コンピュータシステム及び関連するソフトウェアを含んでもよく、プロセッサ、システムメモリ、ストレージデバイス、ネットワークインタフェース、又は他の適切なコンポーネントなどの要素を含んでもよい。例えば、ホスト 400 としての役割を果たすように構成されるコンピュータシステム、又はノードの一実施形態は、図 29 の説明とともに論考される。

#### 【0117】

一般的に、各キーマップインスタンス 140 は、ストレージサービスシステム内に保存されるすべてのオブジェクト 30 について、キーマップ階層を指数化し、管理するために使用されるあらゆる他のデータのみならず、キーマップエントリ 144 を含む、キーマップデータの完全な表現を維持するように設定されてもよい。キーマップインスタンス 140 内で、キーマップデータはホスト 400 にわたって分散されてもよく、個々のホスト 400 はキーマップデータの一部の部分（おそらく冗長）を保存してもよい。図 11A に、ただわずかのホスト 400 が示されるが、他の実施形態において、各キーマップインスタンス 140 は、あらゆる適切な数のホスト 140 を有してもよいことに留意されたい。例えば、一部の大規模な実装において、数ダース、又はことによると数 100 のホスト 140 が、キーマップインスタンス 140 に含まれてもよい。一部の実施形態において、所定のキーマップインスタンス 140 のためのホスト 400 は、所定の領域 310、又はデータセンタ 300 内にローカライズされるが、他の実施形態において、このようなホスト 400 は、異なる領域 310、又はデータセンタ 300 間に分散されてもよいことを考慮されたい。さらに、ホスト 400 が、一部の実施形態においてキーマップ関連機能性のみを実装するように構成されるが、他の実施形態において、ホスト 400 は、ストレージサービスシステムの他の要素に関連する機能性を実装してもよい。例えば、一実施形態において、様々なホスト 400 の 1 つは、ビットストアノード 160 として設定されてもよく、したがって、オブジェクトデータのみならず、キーマップデータも保存してもよい。

#### 【0118】

図 11B は、さらに詳しくキーマップインスタンス 140 a の例示的な実施形態を示す。図示した実施形態において、キーマップインスタンス 140 a 内の各ノード 400 a ~ c は、それぞれのパーティションインデックス 410 a ~ c 及び任意の数のブリック 415 を含む。一般的に言えば、ブリック 415 は、キーマップインスタンス 140 内の中間キーマップデータ構造と対応してもよい。一部の実施形態において、図 12 の説明とともに以下にさらに詳しく述べるように、キーマップデータは、ブリック 415 の間でパーティションに分割されてもよく、キーマップインスタンス 140 内のパーティションの複製が、ブリックレベルで起こってもよい。パーティションインデックス 410 は、キーマップ操作の間に 1 つ以上の特定のブリックの処理の選択を容易にするために、ブリック 415 を指数化方式にするように設定されてもよい。例えば、パーティションインデックス 410 は、ツリー、又は別の適切なデータ構造として設定されてもよい。一実施形態において、キーマップインスタンス 140 内のさらに深い指数レベルのみならず、パーティションインデックス 410 は、次に来る項に詳しく説明する、層別不平衡ツリー、又はトライ構造として言及するデータ向上の特定の種類の一部分として設定されてもよい。図示した実施形態において、キーマップインスタンス 140 は、さらにキーマップコーディネータ 412 を含む。一般的に言えば、キーマップコーディネータ 412 は、さらに詳しく以下に説明する、キーマップアクセス管理、コンテンツ管理、及び同期化方法、又はプロトコルなどを実装するように設定されてもよい。キーマップコーディネータ 412 は、ホスト 400 から区別して説明されるが、一部の実施形態において、1 つ以上のホスト 400 内の処理、又は各要素として実装される場合があることに留意されたい。一部の実施形態において、パーティションインデックス 410 は、ホスト 400 内に別々に実装されるよりはむしろ、キーマップコーディネータ 412 内に実装される場合があることも留意されたい。

10

20

30

40

50

## 【0119】

図11Cは、ブリック415a~nを含む、ホスト400aの例示的な実施形態を示す。示すように、各ブリック415a~nは、任意の数のブロック425のみならず、それぞれのブロックインデックス420a~nを含む。一般的に言えば、ブロック425は、キーマップインスタンス140内で、ブリック415と類似するが、抽象概念のブリックレベルに対して下位である、中間キーマップデータ構造と対応してもよい。パーティションインデックス410と類似して、ブロックインデックス420は、ブリック415内のブロック425の指数化のために構成されるあらゆる適切なデータ構造であってもよい。例えば、ブロックインデックス420は、一実施形態において、層別不平衡ツリーの部分として設定されてもよい。

10

## 【0120】

図11Dに示すように、一実施形態において、ブロック425は、選択のためにインデックスエントリ144を構成するエントリインデックス430のみならず、個々のキーマップエントリ144a~nの任意の数を含むように設定されてもよい。前に説明したように、各エントリ144a~nは、それぞれの記録148a~nのみならず、それぞれのキー146a~nの表示を含んでもよい。

## 【0121】

図11A~Dに示す実施形態のキーマップインスタンス140とキーマップエントリ144との間の階層層間における関係は、図12に要約する。複数のキーマップインスタンス140を含む抽象概念の展開レベルにおいて、特定のキーマップインスタンス140は、抽象概念のインスタンスレベルでパーティションインデックス410を参照してもよい。参照されるパーティションインデックス410は、特定のエントリ144と対応する1つのブリック、又は複数のブリック415を識別してもよい。例えば、図示した実施形態において、すべてのキーマップエントリは、相異なるブリック415と対応する3つの相異なるパーティションによって複製されてもよい。入れ替わりに、所定のブリックは、ブロックインデックス420を介し、特定のブロック425（図12にイは示されない）を参照してもよく、参照されるブロックは、エントリインデックス430を介して特定のエントリ144を参照してもよい。キーマップが、図12に示されるように階層の実装を使用して実装されるが、他の実装も可能であることに留意されたい。大まかに言えば、キーマップインスタンス140は、記録148を有するキー144と関連するあらゆる適切な技術を使用して実装されてもよい。例えば、一実施形態において、キーマップインスタンス140は、従来のデータ系、又は他の種類の構造化指数を使用して実装されてもよい。

20

30

## 【0122】

図12の実施形態における階層層のいくつかは、他の層が拡張性を提供する様に設定されてもよいが、冗長性を提供するように設定されてもよい（例えば、パーティションレベルでブリック415の複製のみならず、展開レベル内でキーマップインスタンス140の複製）。例えば、複数の相異なるレベルにわたる指数化の分散は（例えば、パーティションインデックス410、ブロックインデックス420、及びエントリインデックス430）、キーマップ展開内で指数化されるエントリが増加する場合、インデックスの各部分が管理可能な方法で増えることを可能にすることで、データ構造の拡大を容易にしてもよい。他の実施形態において、異なる冗長、及び非冗長レベルの組み合わせのみならず、さらなる、又はより少ない階層が用いられてもよいことに留意されたい。

40

## 【0123】

オブジェクト30のように、キーマップ階層の層内の複製の使用は、個々の複製の損失に対する敏感性を減少すること初期設定トレランスを改善してもよい。しかしながら、修正が発生した場合、キーマップデータの複製を同期化する試みがなされない場合、キーマップの正しい（例えば、最も現在の）状態はあいまいになる場合があり、代わりに予測不可能な、又は誤ったシステム操作につながる場合がある。一部の実施形態において、キーマップデータの複製された部分は、更新が、各複製に関して永久に検証可能な形で完成されるまで、キーマップクライアントに完了が報告されなくてもよい、アトミック、又はト

50

ランザクションセマンティクスを使用する厳密に同期の方法で更新されてもよい。アトミック更新セマンティクスが、キーマップデータを非一貫性の状態に更新する可能性を最小限にする、またさらには排除してもよいが、アトミック更新の操作は、かなりの規模の分散された環境で、大幅に機能性が低下する場合がある。例えば、キーマップデータの複製が広範囲に分散された場合、更新操作を完了するために必要とされる全体的な時間を決定付ける最も低速の複製では、クライアントから見た複製アクセスの遅延は大幅に異なる場合がある。さらに、1つの複製が失敗した場合、厳密なアトミック更新セマンティクスは、障害が訂正されるまでクライアントを引き止めることを生じる場合があり、クライアントに対する受け入れ難い遅延につながる場合がある。

#### 【0124】

アトミックプロトコルよりも、より良いクライアント操作を提供し得る他の種類の同期プロトコルは、キーマップ階層内に用いられてもよい。一部の実施形態において、特定のキーマップインスタンス140内の複製に関して一部の種の同期プロトコルが採用されてもよい(例えば、図12に示すような、特定のレベルにおける複製)、ハイブリッド同期の方法が実装されてもよいが、別の種類のプロトコルが、キーマップ展開内において異なるキーマップインスタンス140を同期化するために採用されてもよい。このようなハイブリッド方法は、キーマップ階層内の異なるレベルの複製の使用動学に対して、同期の諸経費を、さらに明確に調整させることを可能にしてもよい。

#### 【0125】

例えば、キーマップデータアクセスは、特定のエントリ144に対して反復される要求は、別のキーマップインスタンス140に対するよりはむしろ、特定のキーマップインスタンス140(例えば、地理的、ネットワークトポロジー、又は別の適切な基準において、要求するクライアントに最も近いインスタンス)に向けられる、参照の局所性を示してもよい。すなわち、所定のキーマップインスタンス140内のキーマップデータの複製が、異なるキーマップインスタンス140における対応するキーマップデータよりはむしろ、所定のクライアントによってアクセスされてもよいということになる場合がある。同様に、一部の実施形態において、所定のキーマップインスタンス140内の複製は、相異なるキーマップインスタンス140を同期化するために使用されるプロトコルよりもさらに迅速に収束するように(例えば、複製間の変更を広める)設定されてもよいプロトコルを使用して、同期化されてもよい。

#### 【0126】

一実施形態において、所定のキーマップインスタンス140内のキーマップデータ複製の同期化は、適切なバージョンの定数プロトコルを使用して行われてもよい。一般的に言えば、定数プロトコルに従い操作されるキーマップデータの複製の更新、又は修正(キーマップエントリ格納及び削除操作を含む)は、少なくとも定数の数の複製に関して修正が永久的に(例えば、完全に、及び持続的に)行われた時に、要求するクライアントに関して完了すると見なされてもよい。同様に、定数プロトコルに従い行われるキーマップエントリ取得操作は、同一のデータが、少なくとも定数の数の複製から読み込まれた時に完了すると見なされてもよい。一部の実施形態において、定数の数は、存在する複製の数の単純過半数として定義されてもよいが、他の実施形態において、任意の程度の圧倒的多数が用いられてもよい。定数プロトコル操作は、定数要求が満たされない場合は完了しない場合があることに留意されたい。しかしながら、複製の定数の数が複製の合計数よりも少ない場合、所定の定数プロトコル操作が失敗する可能性は、定数よりも複製間での一致を事実上必要とするアトミックプロトコル操作よりも低い場合がある。本願に説明されるもの以外の定数プロトコルは、キーマップインスタンス140によって採用されてもよい。例えば、Paxosなどの複数相コミットプロトコル、又は2相コミットは、定数種キーマップのセマンティクスに用いられる場合があることに留意されたい。

#### 【0127】

定数プロトコルに従った読み込み及び更新操作の通常の手順において、例えば、通信障害、又は複製の根本をなす供給源の障害により、各複製に更新が伝えることに失敗す

10

20

30

40

50



る可能性がある。一実施形態において、複製間における不一致は、読み込み操作の間に検出され、修復されてもよい。とりわけ、キーマップエントリ取得操作の間に、特定のエントリ144の異なる複製間で異なる値が検出された場合、キーマップ格納操作は、違いを調整するために作成されてもよい。一実施形態において、格納操作の基準として使用されるエントリ144は、異なる値読み込み間におけるタイムスタンプと関連する最新の（例えば、数的に、又は辞書編集的に最も高い）エントリであってもよい。したがって、複製間の相違は、相違の修復のための相異なる過程、又は操作を必要とすることなく、例えば、キーマップエントリ取得操作が処理される場合、「オンザフライ」で解決されてもよい。

#### 【0128】

定数プロトコルの実装を構成するキーマップインスタンス140の実施形態に関するキーマップエントリ格納、取得、削除、及びリスト操作の例示的な実施形態の操作は、図13から図14に示される。様々な実施形態において、これらの方法は、例えば、キーマップインスタンス140内に含まれる1つ以上のホスト400内に設定されてもよい、又は図11Bに示すキーマップコーディネータ412のような、キーマップインスタンス140内の別個の過程、又はシステムとして構成される、キーマップコーディネータ処理内に実装されてもよい。最初から図13までに言及して、キーマップエントリ格納操作は、該操作がコーディネータ120又は他のキーマップクライアントからキーマップインスタンス140で受信される時、ブロック1300で始まってもよい。例えば、特定のオブジェクト30の対応するオブジェクトインスタンスを特定のビットストアノード160に保存することに  
20  
20

#### 【0129】

キーマップインスタンス140の階層は、キーマップエントリ格納操作と対応する複製を識別するために誘導されてもよい（ブロック1302）。例えば、図12の実施形態について、パーティションインデックス410は、どのブリック415が、重要なオブジェクト30と対応するエントリ144を複製するかを決定するように設定されてもよい。続いて、個々の格納操作が、識別された複製に指示されてもよい（ブロック1304）。各格納操作について、キーマップインスタンス140の残存する階層は、対応するエントリ1  
30  
30

#### 【0130】

個々の複製格納操作の成功表示は、複製の定数の数が無事に更新されたかを決定するために監視されてもよい（ブロック1310）。例えば、3つの複製を含む実施形態において、キーマップエントリ格納操作の完了のための複製の定数の数は、2であってよい。複製の定数の数が無事に更新された場合、要求されたキーマップエントリ格納操作が完了した表示は、要求するクライアントに返されてよい（ブロック1312）。そうでない場合、監視が続いてもよい。一部の実施形態において、キーマップエントリ格納操作が、処理の開始後特定の期間内に完了しない場合、操作は終了され、エラー表示が要求するクライアントに返される、タイムアウトが強制されてもよい。他の実施形態において、キーマップエントリ格納操作は、それが完了するまで永久に保留を維持してもよい。

#### 【0131】

一実施形態において、キーマップエントリ削除操作は、格納操作の特別な場合として実装されてもよい。このような実施形態においてキーマップエントリ144は、削除センチネル又はフラッグフィールドを構成する付加的なフィールドを含んでもよく、削除操作は、削除  
50  
50

フィールドをアサートされた状態（例えば、「1」などの特定の値をフィールドに設定することによって）に設定するように構成される、格納操作として実行されてもよい。アサートされた削除フィールドを有するこれらのエントリ144は、将来のキーマップ操作中に無視されてもよい。このような一部の実施形態において、アサートされた削除フィールドを有するこれらのエントリ144を消去するためにキーマップインスタンス144を通じて独立して反復されるように設定されてもよい。他の実施形態において、当該エントリ144は、過去のキーマップ動作のログとして永久に保持されてもよい。

#### 【0132】

キーマップエントリ格納操作の方法の一実施形態は、図14に示す。操作は、取得操作が、コーディネータ120又は他のキーマップクライアントからキーマップインスタンス140で受信される時、ブロック1400で始まってよい。例えば、特定のキーと対応するオブジェクトデータに対するストレージクライアント50からの要求に応じて、ノードピッカー130又はコーディネータ120は、特定のキーと対応するロケータを取得するためにキーマップエントリ取得操作を作成し、それによってビットストアノード160が、前項に説明されるようにオブジェクトデータを取り出すためにアクセスされてもよい。

10

#### 【0133】

キーマップエントリ格納操作のように、キーマップインスタンス140の階層は、キーマップエントリ取得操作と対応する複製を識別するために誘導されてもよい（ブロック1402）。続いて、個々の取得操作が、識別された複製に指示されてもよい（ブロック1404）。各取得操作について、キーマップインスタンス140の残存する階層は、対応するエントリ144へのアクセス及び取り出しを誘導されてもよい（ブロック1406）。エントリ144の所定の複製が無事に取り出されると、対応する取得操作は成功を表示してもよい（ブロック1408）。上記、及び図13に示す個々の格納操作のように、個々のエントリ144のそれぞれの複製を目標とする取得操作は、同時に行われてもよく、ブロック1406～1408は、同様に平行で示される。

20

#### 【0134】

個々の複製取得操作の成功表示は、複製の定数の数が無事に読み込まれたかを決定するために監視されてもよい（ブロック1410）。そうでない場合、付加的な複製が読み込まれるまで監視が続いてもよい。キーマップエントリ格納操作について上記に記載されるように、一部の実施形態において、キーマップエントリ取得操作は、定数の数の複製が無事に読み込まれるまで永久に待機してもよい。他の実施形態において、キーマップエントリ取得操作は、一定期間の後、タイムアウトになってもよく、エラー表示後及び/又はその時点で利用可能な最良のデータ（例えば、最も最近のタイムスタンプを有する複製データ）が、要求するクライアントに返されてもよい。

30

#### 【0135】

複製の定数の数が無事に読み込まれた場合、取り出された複製の内容が異なるかどうか決定されてもよい（ブロック1412）。例えば、要求されたエントリ144の各複製の全部は、他の取り出された複製のそれぞれと比較されるか、又はエントリ144のあるフィールドのみ（例えば、記録148のあるフィールド）比較されてもよい。該比較に使用された基準に従い、取り出された複製間に相違がない場合、取り出されたデータは、キーマップエントリ取得操作が完了したことの表示とともに、要求するクライアントに返されてもよい（ブロック1414）。

40

#### 【0136】

複製間の相違が存在する場合、選択基準に従い、複製の1つが選択されてもよい（ブロック1416）。例えば、該基準は、最も高いタイムスタンプ値を有する複製が選択される、各複製のタイムスタンプ値を考慮することを含んでもよい。キーマップエントリ格納操作は、選択された複製のデータを使用して開始されてもよい（ブロック1418）。例えば、格納操作は、上記に記載されるように図13に従い行われてもよい。格納操作の結果として、初めに要求されたエントリ144の複製の定数の数は、選択された複製の内容

50

で書き込まれてもよく、将来の取得操作が、複製間の不一致に遭遇するのをを減少する。格納操作に続き、選択された複製のデータが、キーマップ取得操作が完了した表示と共に、要求するクライアントに返されてよい（ブロック 1414）。一部の実施形態において、複製間で検出された不一致の場合の取得操作の完了は、不一致を解決するための開始された格納操作の完了を条件としてもよいが、他の実施形態において、取得操作は、カットして生じる格納操作が完了したかどうかに関係なく、完了したことを要求するクライアントに表示する

#### 【0137】

上記に記載されるように、一部の実施形態において、キーマップAPIは、検索パターンなど、一部の基準を満たすキーマップエントリ144のこれらのキー146を示すように構成される、キーマップエントリリスト又はカウント操作をサポートしてもよい。一実施形態において、リスト及び/又はカウント操作は、キーマップエントリ取得操作の特別な場合として実装されてもよく、所定のリスト、又はカウント操作の基準を満たす各エントリ144について、対応するキーマップエントリ取得操作が行われてもよい。しかしながら、定数プロトコルに従い、複数の複製からエントリデータ（例えば、記録148）を実際に取り出すことの付加的な諸経費は、キーマップエントリリスト、又はカウント操作に不必要である場合がある。したがって、一部の実施形態において、定数プロトコルに係るキーマップエントリ取得操作のこれらのステップは、キーマップエントリリスト、又はカウント操作から省略されてもよい。例えば、所定のエントリのすべての複製を識別し、ブロック1402～1404の各複製のための個々の取得操作を作成するよりはむしろ、リスト、又はカウント操作について、単一の複製は（例えばブロック415）任意の選択されてもよく、その対応する階層が、リスト、又はカウント操作基準を満たす各エントリ144を識別するために誘導されてもよい。該基準を満たす結果として得られるエントリ144について、結果として得られるエントリ144の対応するキー146、又はカウントは、図14の定数関係の処理部分を回避して、要求するクライアントに返されてもよい（例えば、ブロック1410～1418）。

#### 【0138】

一部の実施形態において、キーマップインスタンス140は、インデックスエントリ144に対して使用された様々なデータ構造に加え、キャッシュを実装してもよい。例えば、キャッシュは、頻繁に使用されるエントリ144のキーに命令されたキーマップ操作が、対応するエントリ144に直接アクセスするために、インデックスデータ構造の誘導を回避することを可能にしてもよく、キーマップエントリ取得操作の性能を改善してもよい。さらに、キャッシュは、普及した、頻繁にアクセスされるキーに関連するホスト400を、キーマップ要求トラフィックによって過負荷になることから防ぐ助けとなってもよい。例えば、キーマップキャッシュがホスト400の間に分散される一実施形態において、キーのコピーは、キーのためのインデックスデータ構造を維持するホスト400よりも異なるホストでキャッシュされてもよい。ホスト400の間でキャッシュするキーのこのような分散に関わらず、キー処理作業負荷は、ホスト400の間でさらに均一に分けられてよい。

#### 【0139】

一実施形態において、キーマップキャッシュは、キーそのものよりも、キー148のハッシュによって保存及びインデックス化されるように設定されてもよい。データハッシュは、不均衡インデックスデータ構造の論考と共にさらに以下に詳しく説明するように、キーマップキャッシュ内で管理することがさらに容易であるような、キー148などの、確定された長さのデータ構造において、様々な長さのデータを表すための効果的な技術を構成してもよい。さらに、様々なハッシュアルゴリズムは、当初は均一に分散されていない場合があるデータ（例えば、共通するデータの多量の部分を有するキー148の一式）のための均一に分散されたハッシュ値を作成してもよく、ホスト400間におけるキーマップキャッシュデータの一定の分散を容易にしてもよい。一部の実施形態において、エントリ144の内容は、対応するキー148のハッシュされた値と共にキーマップキャッシュ

に保存されてもよい。他の実施形態において、エントリ 1 4 4 それ自体の内容よりはむしろ、エントリ 1 4 4 のためのポインタ、又は他の参照情報が保存されてもよい。

【 0 1 4 0 】

一般的に言えば、キーマップキャッシュを含むキーマップ実施形態において、キーマップエントリ格納及び取得操作は、上記説明に対するわずかな修正をとまって行われてもよい。一実施形態において、キーマップエントリ取得操作は、取得操作が、キャッシュのデータ常駐から情報提供されることができかどうかを決定するために、最初にキャッシュを参考にしてもよい。所得操作は、読み込みのための定数プロトコルに移る前に、キャッシュからの応答のために一定量の時間待機してもよい。キャッシュが、定数プロトコル読み込みが開始された後に値を返した場合、キャッシュからの値読み込みが処理され、対応するエントリ 1 4 4 が返され、定数プロトコル読み込みが終了されてもよい。キャッシュから値が返されない場合、定数プロトコル読み込み操作からのエントリ 1 4 4 読み込み、又は当該のエントリ 1 4 4 に対するポインタは、対応するキー情報と共にキーマップキャッシュにインストールされてもよい。

10

【 0 1 4 1 】

一般的に言えば、キャッシュを含むキーマップ実施形態におけるエントリ格納操作は、同一のキャッシュエントリを修正するために、複数の格納操作を同時に試みることを妨げるためにロックされ、又は他の一貫性プロトコルが採用されることを除き、上記に記載されるように実質的に行われてもよい。一実施形態において、キーマップエントリ操作は、書き込みのための定数プロトコルを開始する前に、キー 1 4 8 と対応するキャッシュエントリをロックすることを試みるように設定されてもよい。ロック要求が成功したキャッシュからの応答を受け取ると（例えば、エントリに他にロックが存在しないため、又はキャッシュに対応するエントリがないため）、定数プロトコルを続行してもよい。定数プロトコルに従い納操作が完了した後、ロックが解放され、新規エントリデータがキャッシュにインストールされてもよい。

20

【 0 1 4 2 】

一部の実施形態において、今しがた説明したように、キーマップエントリ格納及び取得操作のための定数プロトコルは、キーマップエントリ状態を更新するための強一貫性モデルを実装してもよい。すなわち、定数プロトコルは、特定のキーに対する格納操作が、完了するときにクライアントに認識されるとそれに続く取得操作が、取得操作が続行された時点において、各複製が更新されてない場合でも、最も最近格納されたデータを返すことを保証してもよい。

30

【 0 1 4 3 】

格納及び削除操作などの、キーマップ操作は、特定のキーマップインスタンス 1 4 0 にたいして命令され、その特定のキーマップインスタンス 1 4 0 内のエントリ 1 4 4 の状態は、時間と共に変わってもよい。したがって、それらを調整する試みがないと、展開内の異なるキーマップインスタンス 1 4 0 は、時間と共に相違、又は不一致になる傾向がある場合がある。たった 1 つのストレージサービスクライアント 5 0 が、所定のオブジェクト 3 0 を参照する場合、またそれを同一のキーマップインスタンス 1 4 0 を介して行う場合、当該の相違は実用上の影響を有さない場合がある。しかしながら、複数のストレージサービスクライアント 5 0 が、異なるキーマップインスタンス 1 4 0 を介して同一のキーを参照する場合、当該の不一致は、クライアント 5 0 が、同一の時点において、異なるキーマップ状態及び/又はオブジェクトデータの異なるバージョンを観察する結果となる場合がある。

40

【 0 1 4 4 】

前に説明したように、アトミックプロトコル、又は定数プロトコルのような強一貫性プロトコルは、クライアントが複製の不一致を観察することを効果的に防ぐために、又は当該の不一致が起こることを完全に防ぐために、複製を更新する時に用いられることがある。しかしながら、異なる複製のアクセス待ち時間が異なる分散された状況において、時々多量の強一貫性プロトコルが非常に高い操作コストを有してもよい。例えば、アトミック

50

プロトコル、又は定数プロトコルについて、操作完了に必要なとされる時間は、それぞれ、すべての複製の最も遅いもの、又は複製の定数の数の最も遅いものに関して、操作を完了するために必要とされる時間の機能であってもよい。さらに、強一貫性プロトコルがないことにおいて、クライアントに対して複製の不一致が可視になることの見込みは（例えば、ストレージサービスクライアント 50 が古くなったキーマップ、又はオブジェクトデータを得ることの可能性）、概して、アクセスされた複製がまだ更新を反映しない期間内に、クライアントが複製にアクセスすることの見込み機能であってもよい。

#### 【0145】

多くのオブジェクト 30 について、この後者の見込みは低くてもよい。例えば、一部の例において、大部分のオブジェクト 30 は、特定のキーマップインスタンス 140 を介する単一のクライアント 50 によってアクセスされてもよいストレージサービスシステムによって管理され、この場合、不一致は、クライアントの観点からは現実的な価値を無くす場合がある。複数のクライアント 50 からアクセスされてもよいオブジェクト 50 について、観察できる不一致は、依然として可能性が低い。例えば、2つのキーマップインスタンス 140 は、例えば、10秒間の期間、特定のキーに関して不一致である。しかしながら、不一致の期間に関してアクセスが行われない場合（例えば、対応するオブジェクト 30 のアクセス間の持続時間が不一致の期間よりも長い場合）、又は行われるアクセスがさらに最近更新されたキーマップインスタンス 140 に命令された場合（例えば、キーの状態を最後に更新したクライアント 50 が、同一のキーマップインスタンス 140 を介するキーを参照するための次である場合）、不一致はクライアント 50 に目立った影響を有さない場合がある。したがって、一部の実施形態において、キーマップインスタンス 140 は、キーマップインスタンス 140 を一貫した状態に収束するよう努める緩やかな同期プロトコルを採用してもよいが、いかなる時において、キーマップインスタンス 140 の間にある程度の不一致を与える場合がある。このような同期プロトコルは、厳密な同期化が必要でないような、大部分のクライアント 50 に対して、よりよい全体的な機能を提供してもよい。一部の実施形態において、共有されるオブジェクト 30 のためのキーマップデータの厳密なアクセス同期化を必要とするクライアント 50 は、すべてのクライアント 50 が厳密な同期化の義務を課すことを要求することなく、クライアントの間で付加的なプロトコルを実装してもよい。例えば、特定のオブジェクト 30 の 1組へのアクセスを共有するクライアント 50 の 1組は、キーマップデータへのアクセスを調整するセマフォ、又は他の分散されたロック技術を採用してもよい。

#### 【0146】

一部の実施形態において、キーマップインスタンス 140 の間の緩やかな同期プロトコルは、同期化過程の異なる態様を独立して実行する、異なる同期化タスクの組み合わせを含んでもよい。図 15A ~ B は、2つの相異なる同期化タスク、図 15A に示す更新伝播タスクと、アンチエントロピー、又は図 15B に示す設定調整タスクを含む、緩やかな同期プロトコルの操作の方法の一実施形態を示す。図 15A に最初に言及して、1つのキーマップインスタンス 140 に更新するブロック 1500 で始まる操作が検出されてもよい。例えば、キーマップインスタンス 140 は、上記に記載されるように、定数プロトコルに従い、キーマップエントリ格納、又は削除操作を受信し、完了してもよい。

#### 【0147】

キーマップ更新を処理したキーマップインスタンス 140 は、ストレージサービスシステム内でプロビジョニングされたそれぞれの他のキーマップインスタンス 140 に対する更新操作を転送してもよい（ブロック 1504）。例えば、キーマップインスタンス 140 a がキーマップエントリ格納操作を処理した場合、キーマップインスタンス 140 b と 140 c に引数、パラメータなどを含む操作を転送してもよい。一実施形態において、転送は照合、又は認識なしで行われてもよい。例えば、キーマップ更新操作を処理したキーマップインスタンスは、転送された操作がその送り先で受信されかどうかを照合すること、又はそれが受信されていない場合には操作を再送信することを試みることで互いのキーマップインスタンスに操作を転送することを一度だけ試みる、「ファイアー・アン

10

20

30

40

50

ド・フォーゲット」プロトコルを使用する操作を転送してもよい。このような転送は、源を発するキーマップインスタンス 140 から複数のキーマップインスタンス 140 への同時送信、源を発するキーマップインスタンス 140 から他のインスタンスへの配列転送、ツリー系方策などの、あらゆる適切な転送方策を使用して生じてよい。

#### 【0148】

転送された操作を受け取るこれらの関連するホスト 400 は、更新操作をローカルで実行してもよい（ブロック 1506）。例えば、ホスト 400 f が、ホスト 400 a から転送されたキーマップエントリ格納操作を無事に受信した場合、それはあらゆるキーマップクライアントから操作を受信したかのように操作を実行してもよい。ホスト 400 f で格納操作が無事に完了した場合、その結果として、キーマップインスタンス 140 a と 140 b は、格納操作に関して同期化されてもよい。

10

#### 【0149】

一般的に言えば、ホスト 400 間における転送するキーマップ更新操作は、大多数の場合成功する。したがって、このような操作を転送することを伴う諸経費を最小限にすることは、大多数の事例において、キーマップインスタンス 140 間における同期化を達成するために必要とされる時間及び/又は回線容量を削減する。例えば、転送過程から、肯定応答、又は他の種類の照合プロトコル、又はハンドシェイクを消去することは、さらに大きな程度の同期化トラフィックを伴う、キーマップ実装のさらに大きい規模をサポートすることなど、他のユーザのための伝達回線容量を自由にする。多くの例において、キーマップ展開を通じてキーマップ更新を展開するために必要とされる時間は（例えば、所定のキーマップエントリ 144 の複製の不一致の可能性のウィンドウと概して対応してもよい）、ホスト 400 と関連する操作を転送するために必要とされる伝達の遅延、ホスト 400 が転送された操作を適応する為に必要とされる処理の遅延を制限してもよい。しばしば、この合計時間は、秒の順序、又は秒の分数であってもよい。

20

#### 【0150】

しかしながら、一部の例において、ホスト 400 間におけるキーマップ更新操作の転送は失敗する場合がある。例えば、伝達リンクの障害は、別のホストから 1 つのホスト 400 へ達し得ないことを伝える場合があり、又は転送された操作の損失、不完全化、又は輸送中に破損することを生じる場合がある。代わりに、送り先ホスト 400 は、例えば、一時的なハードウェア又はソフトウェアの問題のため、適切に転送された更新操作を受信し、又は正しく処理することに失敗する場合がある。一実施形態のように、転送されたキーマップ更新操作が、目的とされるホスト 400 によって無事に受信され、処理されたことを照合する、又は確かめるために、源を発するホスト 400 の部分になにも試みがなされない場合、個々の操作の転送の失敗は、あるエントリ 144 に関するキーマップインスタンス 140 間における不一致の結果になる場合がある。

30

#### 【0151】

同様に、一実施形態において、キーマップインスタンス 140 間における緩やかな同期プロトコルは、上記に記載されるように、及び図 15B に示すように、アンチエントロピー、又は設定調整タスクを含んでもよい。このタスクは、概して、タスクの操作は、異なるキーマップインスタンス 140 間における相違点を減少し、類似点を増加する役割を果たす場合があり、したがって、適切に同期化されたインスタンスに対する更新伝播の無作為な、又は体系的な障害によってもたらされることがある、キーマップインスタンス 140 間における全体的なエントロピーの減少となる、「アンチエントロピー」タスクといわれることもある。図示した実施形態において、操作はブロック 1510 で始まり、開始しているキーマップインスタンス 140 が、特定のパーティションの調整を行う、別のキーマップインスタンス 140 を無作為に選択し、異なるホスト 400 における

40

#### 【0152】

開始しているキーマップインスタンス 140 は選択されたキーマップインスタンス 140 を有するインスタンス内のパーティションについての情報を交換してもよい（ブロック 1512）。例えば、2 つのキーマップインスタンス 140 内の特定のホスト 400 は、

50

各インスタンス内で維持されるパーティションインデックス 4 1 0 のコピーを交換し、入れ替わりに各インスタンス内のこれらのブリック 4 1 5 を識別するように設定されてもよい。

#### 【 0 1 5 3 】

交換されたパーティション情報に基づき、開始しているキーマップインスタンス 1 4 0 は、2つのインスタンスのパーティション間での一致を識別してもよく（ブロック 1 5 1 4）、選択されたキーマップインスタンス 1 4 0 内の対応するパーティションを有する源を発するキーマップインスタンス 1 4 0 内の各パーティションを調整してもよい（ブロック 1 5 1 6）。例えば、前に説明したように、所定のキーマップインスタンス 1 4 0 内の各パーティションは、多くのブリック 4 1 5 に渡って複製されてもよい。一実施形態において、開始しているキーマップインスタンス 1 4 0 は、パーティション間の相違を調整するために、パーティション内の特定のブリック 4 1 5 が（「リードブリック」といわれることもある）選択されたキーマップインスタンス 1 4 0 内の対応するパーティションの対応するブリック、又は「ピア」ブリック 4 1 5 と通信することを命令するように設定されてもよい。一実施形態において、2つのブリック 4 1 5 の調整は、ブリックが各ブリック 4 1 5 に含まれるキーマップエントリ 1 4 4 の相違についての情報を交換することを伴い、各キーマップインスタンス 1 4 0 内の最も現在の情報を展開することを伴ってもよい。例えば、1つのブリック 4 1 5 が、タイムスタンプ情報に基づき、エントリ 1 4 4 のバージョンがピアブリック 4 1 5 よりもさらに現在のものであることを決定する場合、ブリックはエントリデータからピアブリック 4 1 5 に対して通信してもよい。それに続き、ピアブリック 4 1 5 は、エントリ 1 4 4 のそのコピーを更新するために、キーマップエントリ格納操作（例えば、上記に詳しく説明した定数プロトコルに従い）を行ってもよい。

#### 【 0 1 5 4 】

2つのキーマップインスタンス 1 4 0 間のパーティション調節が完了すると、操作は、調整操作が、別のランダムキーマップインスタンス 1 4 0 に関し再び開始されるブロック 1 5 1 0 から継続されてもよい。様々な実施形態において、各キーマップインスタンス 1 4 0 は、予め所定の、又は動的に決定された間隔でこの課程を実行するように設定されてもよい。例えば、調整は、1分に1回の静的割合、又はランダム確率分布、又は他の統計的確率分布に従い決定される間隔で起こってもよい。一部の実施形態において、ある数のキーマップアクセスが起こった後に、又はある個々の1つにアクセスした後に、キーマップエントリの種類、又はグループが検出された後に調整が行われてもよい。

#### 【 0 1 5 5 】

一般的に言えば、図 1 5 A ~ B に示す更新伝播及び調整の設定又はアンチエントロピーの方法は、補足的な方法で行われてもよい。大部分の状況下において、更新伝播は、展開内において異なるキーマップインスタンス 1 4 0 を満足に同期化してもよい。キーマップの不一致が更新伝播の障害によりもたらされるこれらのインスタンスにおいて、アンチエントロピータスクは、当該の不一致を調整するために概して行われてもよい。一部の実施形態において、アンチエントロピータスクの実行は、2つのキーマップインスタンス 1 4 0 がその全体において正確に同期化されることを保証しない場合があることに留意されたい。しかしながら、一実施形態において、アンチエントロピータスクは、その実行が2つのキーマップインスタンス 1 4 0 間の不一致の程度を増加しないことを保証するために実装されてもよい。したがって、度重なるアプリケーションの全体に渡って、アンチエントロピータスクは、キーマップインスタンス 1 4 0 の収束を容易にしてもよい。アンチエントロピータスクの一実施形態におけるさらなる詳細は、キーマップインスタンス 1 4 0 が実装されてもよいデータ構造の特定の実施形態の説明とともに以下に提供される。

#### 【 0 1 5 6 】

図 2 に示し、上記で論考したように、一部の実施形態において、ストレージサービスシステムは、他のキーマップインスタンス 1 4 0 に加え、レプリケータキーマップインスタンス 1 9 0 を含んでもよい。一実施形態において、レプリケータキーマップインスタンス 1 9 0 は、上記に記載されるとおり、キーマップインスタンス 1 4 0 と基本的にあらゆる

点で等しく構成され、上記で論考したように、プロトコルを使用してキーマップ同期化に  
関与してもよい。しかしながら、このような実施形態において、レプリケータインスタンス  
190は、コーディネータ120、又は他のキーマップクライアントよりはむしろレプリ  
ケータ180の役目をするように設定されてもよい。一部の状況において、他のキーマッ  
プインスタンス140から分離するレプリケータキーマップインスタンス190は、概し  
てキーマップ能力を改善してもよい。例えば、レプリケータ180は、オブジェクト30  
の複製のヘルス状態及び数をチェックするためにキーマップを通じて反復する、相当量の  
キーマップ要求トラフィックを作成してもよい。ストレージサービスクライアント50の  
要求を代行して作成されたキーマップトラフィックと混合されると、レプリケータキーマ  
ップトラフィックは、応答時間、又は他のクライアント50に関するサービスの質の目安  
にマイナスの影響を与える場合がある。それに反して、戦況のキーマップインスタンス1  
90を使用するためにレプリケータ180を構成することは、クライアント作成のトラフ  
フィックから内部で作成されたキーマップトラフィックを分離させてもよい。加えて、この  
ような分離は、その大部分のクライアントの要求に従い、各種類のキーマップインスタン  
スを増やすことを実装することをより可能にする。例えば、レプリケータキーマップイン  
スタンス190の実装は、いかなる所定のキーマップ操作の遅延を最小にするためよりは  
むしろ、多数の同時キーマップ操作の処理を容易にするように設定され、キーマップイン  
スタンス140は、サービスの質基準の異なる組み合わせのために最適化される。しかし  
ながら、この方法におけるキーマップインスタンスの分離は必要ではなく、一部の実施形  
態対において、レプリケータ180は、専用のレプリケータキーマップインスタンス19  
0よりはむしろ、キーマップインスタンス140のクライアントであってもよいことに留  
意されたい。

#### 【0157】

一実施形態において、レプリケータキーマップインスタンス190は、クライアント5  
0によるストレージサービスシステム供給源の使用の会計を容易にするように設定されて  
もよい。具体的に、レプリケータキーマップインスタンス190は、対応するオブジェク  
ト30に対する請求、又は他の金融責任を負うそれぞれのエンティティを示す付加的なデ  
ータを有する、キーマップインスタンス140によって保存されたエントリ144を増や  
すように構成されてもよい。例えば、図16に示した実施形態において、レプリケータキ  
ーマップエントリ194が示される。entry 194 may function identically to entries  
144 with respect to the structure and hierarchy of keymap instances 140レプリケ  
ータキーマップインスタンス190内において、エントリ194は、キーマップインスタン  
スの構造及び階層に関して、エントリ144に対してあらゆる点で等しく機能してもよ  
い。しかしながら、具体的な実施形態において、エントリ194は、付加的なフィールド  
、バケットID196を含む。一般的に言えば、バケットID196は、キー146に対  
応するオブジェクト30を含むバケット20の識別子の表示を含んでもよい。当該識別子  
は、例えば、上記に記載されるようにオブジェクト30を保存するためのバケット20を  
作成するためにクライアント50からの要求に応じる、ウェブサービスインターフェース  
100又はコーディネータ120によって定義されてもよい。他の実施形態において、会  
計情報は、レプリケータキーマップインスタンス190のエントリ内だけを反映する必要  
はないことに留意されたい。例えば、一実施形態において、一部の、又はすべてのキーマ  
ップインスタンス140のキーマップエントリ144は、例えば、記録148、又はキー  
146内の付加的なフィールドとして、バケットID196の表示を保存するように設定  
されてもよい。

#### 【0158】

上記に記載されるように、オブジェクト30とバケット20との間の関係は、キーマッ  
プインスタンス140の全体的な操作に対し、透過的であってもよい。しかしながら、も  
し、この関係が一般的に静的である場合、レプリケータキーマップエントリ194を介し  
て明示的に関連するバケット20及びオブジェクト30は、クライアント50の会計及び  
請求を容易にする場合がある。例えば、各オブジェクト30と関連するバケット20のた

10

20

30

40

50



めのウェブサービスインターフェース 100 を明示的に要求するよりはむしろ、会計過程（レプリケータ 180、又は別のモジュールに含まれてもよい、又はシステム内の相異なるボジュール内に実装されてもよい）は、バケット ID 196 に従い、レプリケータキーマップエントリ 194 をソートするように設定されてもよい。当該のソートが完了すると、特定のバケット ID と関連するすべてのキー 146 は容易に明らかになる。記録 148 内に表示されたように対応するオブジェクト 30 のサイズは、バケット ID 196 に関連する全体のストレージリソースの使用を決定するために合計される。さらに、特定のオブジェクト 30 と関連するストレージのクラスなど、オブジェクト 30 の他の特徴が配慮されてもよい。供給源の使用が、適切な請求モデルに従い通貨に定められてもよい。

【0159】

10

様々な実施形態において、レプリケータキーマップエントリ 194 は、様々な内部システム管理、又は会計タスクを容易にしてもよいバケット 196 のかわりに、又はそれに加えて他のフィールドを含んでもよい。レプリケータキーマップインスタンス 190 が他のキーマップインスタンス 140 から相異なる実施形態において、当該付加的なフィールドのストレージコストはレプリケータキーマップインスタンス 190 を構成してもよいことに留意されたい。しかしながら、専用のレプリケータキーマップインスタンス 190 を欠く実施形態において、キーマップインスタンス 140 のエントリ 144 は当該付加的なフィールドを含むために増大されてもよい。

【0160】

#### 層別不平衡データ構造

20

前に説明したように、一部の実施形態において、ストレージサービスシステムは、例えば、数十億又はそれ以上の、非常に多いオブジェクト 30 の数をサポートするために増やしてもよい。したがって、このような実施形態において、各キーマップインスタンス 140 は、管理するための同様の数のエントリ 144 を有するであろう。一部の実施形態において、キーマップインスタンス 140 は、前項で論考したキーマップエントリリスト、及びカウント操作などの、様々な種類のソート化及び/又はグループ化操作をサポートしてもよい。さらに、一貫したキーマップ操作をサポートするために、各キーマップインスタンス 140 によって管理される多くのキーは、上記に記載されるように他のキーマップインスタンス 140 の間で同期化される必要があってもよい。

【0161】

30

多くの状況において、キーマップインスタンス 140 によって提供されたキーマップ機能性は、全体的なストレージサービスシステムの操作の中心である。例えば、クライアント 50 がオブジェクト 30 の特定のインスタンスに対してロケータ系のアクセスを行わないことを選ぶ場合、キーマップインスタンス 140 はクライアント 50 によって操作される各キー系のオブジェクトアクセスを仲介してもよい。したがって、クライアント 50 によって見られるストレージサービスシステムの操作は、アクセスするキーマップインスタンス 140 及びキーマップエントリ 144 の処理の効率及び速度に直接左右される場合がある。言い換えると、キーマップインスタンス 140 の性能は、図 12 の実施形態の、パーティションインデックス 410、ブロックインデックス 420、エントリインデックス 430 を実装するために使用されるデータ構造などの、エントリ 144 をインデックス化し、体系化するデータ構造に直接左右される。

40

【0162】

大規模のキーマップ実装において、ソート化及び同期操作をサポートするためのインデックスデータ構造を設計することは、相当量の課題をもたらす場合がある。例えば、データ系など、大量のデータのインデックス化を必要とする従来のアプリケーションは、B ツリー、又は他の種類の平衡ツリーなどの従来の平衡データ構造をしばしば用いる。一般的に言えば、キーマップエントリ 144 などの、所定の数量のデータアイテムをインデックス化するために使用する場合、平衡データ構造アルゴリズムは、管理すべきアイテムの数量に従い、平衡データ構造にわたって、データアイテムを分散する試みをする。例えば、所定の 10000 のキーマップエントリ 144 をインデックス化するためには、平衡デ

50

ータ構造アルゴリズムは、エントリが1グループおよそ1000のエントリの10のグループに分割されるような、エントリ144の間の区切り点を選ぶ試みをしてよい。平衡データ構造アルゴリズムは、例えば、およそ1000のエントリの各グループを、それぞれおよそ200のエントリの5つの下位グループに再分割する、各グループ内に平衡した階層のさらなるレベルを作ってもよい。データアイテムが平衡データ構造に加えられ、データ構造から削除されると、データ構造内のグループ及び/又は下位グループは不平衡となる場合がある。したがって、従来の平衡データ構造アルゴリズムは、グループ間でデータアイテムを再配置することによって、付加的なグループを作ること、及び/又は階層の付加的なレベルを作ることによって、データ構造を再平衡する。このような再平衡は、データアイテムが加えられ、又は削除された場合、「オンザフライ」が行われる、又はある数のデータアイテム修正が行われたあと、又は最後の再平衡から一定期間が経過した後に「オンザフライ」が起こる場合がある。

10

#### 【0163】

平衡の取れた方法においてデータを分離することの長所によって、平衡データ構造は、予測可能な、データ構造内でのいかなる所定のデータアイテムのおよそ一定のアクセス遅延をもたらしてもよく、多数のデータアイテムをインデックス化することが必要である、大規模の実装において望ましい場合がある。しかしながら、例えば、上記に記載されるように緩やかな同期化モデルを使用することは、平衡データ構造の分散されたインスタンスを効果的に調整、又は同期化することは特に困難である場合がある。具体的には、平衡データ構造のインスタンスが独立して修正された場合、各インスタンス内でデータアイテムをグループに分割する区切り点が不一致になる場合がある。その結果、異なる平衡データ構造インスタンスのグループ又は下位グループ間のデータアイテムメンバーシップに関して、直接の一致がない場合がある。2つの当該インスタンスを調整するためには、2つのインスタンスの全体を徹底的に比較することが必要であることがあり、各インスタンスインデックスが多数のデータアイテムである場合は非常に時間がかかる場合がある。

20

#### 【0164】

数量に従いデータアイテムをグループ間に分散する平衡データ構造の代替として、一部の実施形態において、キーマップインスタンス140のインデックスデータ構造は、各グループ内のデータアイテム間の一部の関係に従い、グループ間でデータアイテムを分散する、非平衡データ構造（トライとしても言及される）を実装するように設定されてもよい。具体的には、キーマップインスタンス140は、対応するキー146のプレフィックスに従い、エントリ144をインデックス化するように設定されてもよい。例として、対応する、大文字と小文字を区別しない英数字キー146を有する600のキーマップエントリ144が存在する事例を考慮されたい。これらの600のエントリの平衡インデックスは、それぞれ200のエントリの3つの平衡グループのエントリに分割されるであろう。それに反して、一実施形態において、非平衡インデックスは、aからlまでの文字で始まるこれらのエントリが、第1のグループに割り当てられ、mからxの文字で始まるこれらのエントリが第2のグループに割り当てられ、y又はzの文字又は数字0～9で始まるこれらのエントリが第3のグループに割り当てられる、3つの英数字グループを定義するであろう。

30

40

#### 【0165】

エントリ144は、非平衡インデックスのグループにわたって不均一に分散されてもよい。例えば、第1のグループに300のエントリ、第2のグループに250のエントリ、そして第3のグループにわずかに50のエントリがある場合がある。しかしながら、いかなる所定のエントリ144について、非平衡の特定のグループにおける所定のエントリ144のメンバーシップは、いかなる特定のグループにおける多くのエントリ144に依存することなく、その対応するキー146の機能であってもよい。したがって、非均衡印デックする2つが、同一のグループ定義を維持する場合、各グループは、一方のグループに依存することなく、独立して同期化されてもよい。例えば、2つのインスタンス間におけるa～lグループは、mからxグループ及びyから9グループから独立して同期化されても

50

よい。それに反して、上記に記載されるように、エントリ 1 4 4 の同一の組の平衡インデックスの 2 つのインスタンスの同期化は、すべてのグループにわたるすべてのエントリを考慮することを必要とする場合がある。

#### 【 0 1 6 6 】

多くのデータアイテムをインデックス化するための非平衡データ構造の使用を示す 1 つの例を、図 1 7 に示す。図示した実施形態において、非平衡インデックス 2 0 0 (又は、単にインデックス 2 0 0) は、プレフィックス「a l」で始まる多くの文字値をインデックス化するための階層的な方法に配置される、多くのノード 2 1 0 a を含む。例えば、インデックス化された値は、キーマップインスタンス 1 4 0 の様々なエントリ 1 4 4 のキー 1 4 6 と対応してもよい。インデックス 2 0 0 内の各ノード 2 1 0 は、インデックス化されるデータアイテムと直接一致してよい、又は一致しなくてよい関連するタグ値を含む。図示した実施形態において、楕円で描かれたノードは、対応するデータアイテムを有しないインデックス 2 0 0 の内部ノードと対応してもよいが、長方形で描かれたノードは、インデックス化されたデータアイテムと対応してもよい。したがって、例えば、ノード 2 1 0 a は文字列「a l」と対応し、インデックス 2 0 0 内の多くの他のノードと関係するが、文字列「a l」と対応する実際のキー 1 4 6 は存在しない場合がある。それに反して、「a l i c i a」というタグを有するノード 2 1 0 n は、同一の文字列を特定するキー 1 4 6 と対応してもよい。内部ノード 2 1 0 と非内部ノード 2 1 0 との間の相違は、ノード 2 1 0 の状態を明示的に反映してよく、又は反映しなくてもよい。

#### 【 0 1 6 7 】

以下に説明するように、一部の実施形態において、非平衡データ構造は、他の複数のインデックスの 1 つのインデックスとして構成されてもよい。このような一部の実施形態において、インデックス 2 0 0 の第 1 のインスタンス内のデータアイテムは、別のインデックス 2 0 0 のルートノード 2 1 0 であってもよく、第 1 のインデックス 2 0 0 内の対応するノード 2 1 0 は、非内部ノードと見なされる場合がある。すなわち、一部の実施形態において、所定のインデックス 2 0 0 の非内部ノード 2 1 0 は、所定のインデックス 2 0 0 の外部である、エントリ 1 4 4、又は別のインデックス 2 0 0 のルートノードなど、データ値と関連するいかなるノード 2 1 0 として概して定義される場合がある。同様に、所定のインデックス 2 0 0 の内部ノードは、所定のインデックス 2 0 0 内の他のノード 2 1 0 のみを参照してもよく、エントリ 1 4 4、又は所定のインデックス 2 0 0 から異なる他の

#### 【 0 1 6 8 】

様々な実施形態において、各ノード 2 1 0 は、様々な情報を暗号化してもよい。ノード内でコード化されてもよい様々なデータフィールドを示す、ジェネリックノード 2 1 0 の一実施形態を、図 1 8 に示す。具体的な実施形態において、ノード 2 1 0 は、タグフィールド 2 1 2、カウントフィールド 2 1 4、フィンガープリントフィールド 2 1 6、及び 1 つ以上のポインタフィールド 2 1 8 を含む。概して、タグ 2 1 2 は、以下にさらに詳しく説明するように、インデックス 2 0 0 をトラバースする又は操る過程で使用されてもよい、所定のノード 2 1 0 と対応する値を保存するように設定されてもよい。一部の実施形態において、タグ 2 1 2 は、インデックス 2 0 0 内のすべてのノード間からノード 2 1 0 を独自に識別してもよい。また、一部の実施形態において、所定のノード 2 1 0 のタグ 2 1 2 は、インデックス 2 0 0 内の所定のノード 2 1 0 のすべての直接の先祖であるタグ 2 1 2 に、プレフィックスとして含んでもよい。すなわち、所定のノード 2 1 0 のタグ 2 1 2 は、その所定のノードの直接の親ノード 2 1 0 のタグに対する一部の値を加えることを通じて決定されてもよい。例えば、タグ、「a l i c i a」を有する、図 1 7 のノード 2 1 0 n を考慮されたい。ノード 2 1 0 n の直接の先祖ノード 2 1 0 l、2 1 0 k、及び 2 1 0 a のそれぞれは、ノード 2 1 0 n のタグの適切なプレフィックスを形成するタグ(それぞれ「a l i c」、「a l i」及び「a l」)を有する。

## 【0169】

図17に示すように、あるノード210は、インデックス200の階層のさらに下方で、1つ以上の子、又は子孫ノード210を参照する。一実施形態において、ポインタフィールド218は、所定のノード210から別のノード210へ、ポインタ、又は参照を反映するデータを保存するように設定されてもよい。例えば、所定のポインタフィールド218は、メモリアドレス空間などの、アドレス空間内で参照されたノード210の位置を識別するアドレスを含んでもよい。また、所定のポインタフィールド218は、参照されたノード210に関する付加的なタグ情報を含んでもよい。例えば、図17に示すように、子孫ノード210への所定のノード210からの各アークは、所定のノード210のタグ212によって形成されるプレフィックスとは異なる子孫ノード210のタグ212の第1の文字をラベルされる。一実施形態において、この付加的なタグ情報は、参照されたノード210に対するポインタと共に、対応するポインタフィールド218内に保存されてもよい。例えば、ノード210aに含まれるポインタフィールド218は、対応するタグデータ、「a」、「e」、「f」、「i」、及び「z」のみならず、ノード210b、210g、210j、210k、及び210tに対する参照をそれぞれ含んでもよい。

10

## 【0170】

図12に関連して上記に記載されるように、インデックス200のようなインデックスは、選択のためにキーマップエントリ144などのデータアイテムを体系化するために使用されてもよい。一部の実施形態において、非内部ノード210のポインタフィールド218（すなわち、インデックス化されたデータアイテムに直接マップするノード210）も、キーマップエントリ144、ブロック425、又はブリック415などの、対応するデータアイテムに対するポインタを含んでもよい。一部の実施形態において、以下にさらに詳しく説明するように、インデックス200のような、非平衡インデックスは、1つのインデックス200の非内部ノード210が、別のインデックス200を参照してもよいような、階層的に実装されてもよい。インデックス化されたデータアイテムを参照するポインタフィールド218は、ポインタフィールド218の異なる種類のために異なるコード化を使用することなどの、あらゆる適切な技術によって、別のノード210を参照するポインタフィールド218から区別されてもよい。例えば、前の段落内で説明したように、子孫ノード210に対するアークに関連するタグ情報がポインタフィールド218内でコード化された実施形態において、空タグは、子孫ノード210への参照から、インデックス化されたデータへの参照を区別するために使用されてもよい。

20

30

## 【0171】

所定のノード210について、カウントフィールド214及びフィンガープリントフィールド216は、所定のノード210の真下のノード210の状態を参照するように設定されてもよい。一実施形態において、カウント214は、所定のノード210の子孫（例えば、階層的に間下にある）であるすべてのノードのカウントを保存するように設定されてもよい。例えば、図17のノード210kは、その真下に、インデックス200内で8つの他のノード210を有する。同様に、そのカウント214は、あらゆる適切なコード化又は形式を使用する、8つの値を表示してもよい。

## 【0172】

様々な実施形態において、所定のノード210のフィンガープリントフィールド216は、所定のノード210の階層的に真下にあるノード210のデータの一部で行われたハッシュの値標示（例えば、適切なハッシュアルゴリズムの結果）を保存するように設定されてもよい。例えば、所定のノード210のフィンガープリントフィールド216は、所定のノード210の子孫であるすべてのノード210のタグ212のハッシュの合計を反映してもよい。代わりに、フィンガープリントフィールド216は、トラバーサル（例えば、横型トラバーサル、又は縦型トラバーサル）の特定の、一貫した順序に従い、子孫ノード210のタグ212の連結のハッシュを反映してもよい。他の実施形態において、タグ212のそばのノード210の他のフィールドは、ハッシュ化に関与してもよい。一部の実施形態において、所定のノード210と関連するデータは、それ自体のフィンガー

40

50

プリントフィールド 216 内で反映されてもよいが、他の実施形態において、所定のノード 210 のフィンガープリントフィールド 216 は、その子孫ノードに基づいて厳密に決定されてもよい。記述の一貫性について、本願で使用した所定のノード 210 のフィンガープリントは、所定のノード 210 の少なくとも一部の子孫ノードの機能であるハッシュ値を参照してもよいが、所定のノード 210 のハッシュは、子孫ではなく、所定のノード 210 のみと関連するデータの機能である、ハッシュ値を参照してもよい。

#### 【0173】

一般的に言えば、ハッシュアルゴリズムは、2つのハッシュ値が異なる場合、2つのハッシュ値が作成された元のソースデータ値も何らかの方法で違わなければならない、小さい、一般的に固定された長さのハッシュ値の上に、あるいは任意の長さの所定のソースデータ値をマップするように設定されてもよい。ハッシュアルゴリズムは、一般的に 1 対 1 の機能ではなく、2つのハッシュ値間のアイデンティティは、元のソースデータ値間のアイデンティティを必ずしも暗示しない。しかしながら、ハッシュアルゴリズムの一部のクラスについて、同一のハッシュ値を与えられた元のソースデータ値間のアイデンティティは、とりわけ、冗長のある程度を示すソースデータ値について、統計的に定量化できる確率、又は信頼度内である可能性がある場合がある。ハッシュアルゴリズムの異なる種類は、署名、フィンガープリント、又はチェックサムアルゴリズムと言われることもある。ハッシュアルゴリズムのいかなる適切な種類は、あらゆる適切なバージョンの Message Digest 5 (MD5) アルゴリズム、又は、SHA-1、SHA-256、SHA-512 などのような Secure Hash Algorithm (SHA) を含む、制限されない例を手段として、フィンガープリントフィールド 216 に保存されるハッシュ値を作成するために用いられてもよいことを意図する。

#### 【0174】

前項に記載のように、キーマップインスタンス 140 で行われる基本的な操作は、操作に対するパラメータとして特定されるキーと対応するエントリ 144 をそれぞれ保存し及び取り出す、格納及び取得操作を含んでもよい。一部の実施形態において、キーマップインスタンス 140 内の様々なインデックスは、インデックス 200 などの非平衡インデックスとして実装されてもよい。

#### 【0175】

多数のデータアイテムをインデックス化する場合、キーマップインスタンス 140 で共通する場合があるように、すべてのデータアイテムのためのインデックス 200 の 1 つのインスタンスを使用することは、実用的でないであろう。例えば、1 つの大きなインデックスは、インデックスを処理するシステムのメモリに完全に収まらない場合があり、インデックスに左右される操作の能力にマイナスに影響を与える場合がある。一部の実施形態において、大きなインデックスは、層別の、非平衡データ構造、又は層別インデックスを使用して実装されてもよい。一般的に言えば、層別インデックス内で、インデックス 200 の複数のインスタンスは、階層的に識別されてもよく、階層が上であるインスタンスは、他のインデックス 200 をインデックス化してもよく、階層が下であるインデックスは、特定のエントリ 144、又は他のエンティティ（例えば、ブロック 425、又はブリック 415）をインデックス化してもよい。

#### 【0176】

層別インデックスの一実施形態を図 19 に示す。図示した実施形態において、層別インデックス 220 は、5 つのインデックス 200 a ~ e を含む。インデックス 200 a は、ノード 210 u ~ x を含み、各ノードはインデックス 200 b ~ e のうちの 1 つのそれぞれのルートノードを参照する、非内部ノードである。代わりに、インデックス 200 b ~ e は、図 17 に示されたノード 210 a ~ t の様々なものをそれぞれ含む。層別インデックス 220 の一部の実施形態において、インデックス 200 a のような、上位インデックスは、メモリ、キャッシュ、又はシステム処理インデックスの別の上位メモリ階層に属するように設定されてもよいが、インデックス 200 b ~ e などの下位インデックスは、ディスク、又は別の当該メモリ階層の下位に主として属してもよい。このような実施形態に

において、下位インデックスは、例えば、ページングタイプの技術を使用して、必要があれば、メモリ階層の下位から上位に再配置されてもよい。多数のデータアイテムのインデックスの階層的なパーティションをサポートすることによって、層別インデックス 220 は、システム供給源をさらに効率的、効果的に使用できる。

#### 【0177】

例えば、前述のページング技術を使用して、層別インデックス 220 の頻繁に使用されたインデックス 200 は、容量が限られるが一般的にアクセスすることが速い上位のメモリ階層に保存されてもよいが、それほど頻繁に使用されないインデックス 200 は、上位よりは大きなストレージ容量を有するが一般的にアクセスすることが遅い下のメモリ階層に保存されてもよい。一部の実施形態において、ノード 210 が、層別インデックス 220 内でインデックス 200 に加えられると、個々のインデックス 200 は、目的とするサイズ（インデックスに実装するシステム上のディスクブロック、又はメモリページなど）を超えて増大する場合があることを考慮されたい。このような実施形態において、所定のインデックス 200 が目的とするサイズを超えて増加した場合、2つ以上のインデックスインスタンスに分割されてもよい。当該の分割の実行過程において、ノード 210 は、新規インデックスインスタンスを説明するために必要であるものとして上位インデックス 200 に加えられてもよい。

#### 【0178】

キーマップエントリの格納又は取得操作に応じて、層別、又は非層別非平衡インデックスは、特定のキーが、インデックス 200 内のノード 210 と対応するかどうかを決定するためにトラバースされてもよい。非平衡インデックストラバーサルの方法の一実施形態を図 20 に示す。図示した実施形態において、インデックス内でキー値を検索する（検索値としても言及される）ブロック 2000 で開始する操作を、例えば、関連するキーマップ操作を介して、特定する。それに続き、インデックスのルートノード 210（例えば、親ノードを有さないノード 210）が選択される（ブロック 2002）。

#### 【0179】

選択されたノード 210 について、ノードの対応するタグ値 212 は、タグ値が、検索値と正確に一致するか、検索値のプレフィックスであるか、又はそのいずれでもないかどうかを決定するために検索値に対して比較される（ブロック 2004）。選択されたノード 210 のタグ値 212 が検索値と一致する場合、選択されたノード 210 は、内部、又は非内部ノードであるかどうか決定するために調査される（ブロック 2006 ~ 2008）。例えば、ポインタ 218 又は選択されたノード 210 の他の内容は、ノードが、エントリ 144、又はインデックス 200 の別のインスタンスなどのインデックス 200 によってインデックス化されたデータ値を参照するかを決定するために調査されてもよい。選択されたノード 210 が内部ノードである場合、以下に説明するようにインデックス不足が生じる場合がある（ブロック 2022）。

#### 【0180】

選択されたノード 210 が非内部ノードである場合、選択されたノード 210 によって参照されたデータ値は取り出される（ブロック 2010）。層別非平衡データ構造をサポートする実施形態において、一部のデータ構造インスタンスは、他のデータ構造インスタンスをインデックス化してもよく、取り出されたデータ値はエントリ 144、又はインデックス 200 の別のインスタンスのルートノードのいずれかと対応してもよい。取り出されたデータ値がエントリ 144 である場合、インデックストラバーサルが完了され、取り出されたエントリ 144 は、トラバースを始めたキーマップ操作に従い、処理されてもよい（ブロック 2012 - 2014）。例えば、開始しているキーマップ操作が取得操作の場合、取り出されたエントリ 144 が取得操作の結果として返されてもよい。例えば、開始しているキーマップ操作が格納操作の場合、取り出されたエントリ 144 が、格納操作で特定されたパラメータに従い修正されてもよい。

#### 【0181】

取り出されたデータ値が、エントリ 144 と対応しない場合、図示した実施形態におい

10

20

30

40

50

て、それは別のインデックス200のルートノード210と対応してもよい。同様にこのルートノード210が選択されてもよく(ブロック2012、2016)、操作は、新規に選択されたインデックス200のトラバーサルであるブロック2004から処理されてもよい。したがって、一実施形態において、図20の方法の実行は、検索値と対応するノード210の存在又は不足が断定的に決定されるまで続いてもよい。

#### 【0182】

ブロック2006に戻って、選択されたノード210のタグ212が検索値と一致しないが、検索値のプレフィックスである場合、選択されたノード210の子孫は、どの子孫も検索値と対応するかを判断するために調査されてもよい(ブロック2018)。そうである場合、対応する子孫ノード210は選択されてもよく(ブロック2020)、操作はブロック2004から処理されてもよい。一実施形態において、選択されたノード210のポインタ218は、選択されたノード210のタグ212と共に取得された時、付加的なタグ情報が特定のポインタ218と関連しているか、また検索値のプレフィックスを形成する(又は完全に一致する)かどうか決定するために調査されてもよい。例えば、図17に言及して、ノード210aのタグ「a1」は、「a1libaba」の検索値のプレフィックスであると決定されてもよい。さらに、対応するポインタ218によって代表されてもよい、ノード210aからノード210kのアーキは、付加的なタグ情報「i」と関連する。このタグ情報は、ノード210aのタグ「a1」に添付されると、検索値のプレフィックスでもある、値「a1i」を形成する。したがって、ノード210kは、将来のトラバーサルのために選択されてもよい。

#### 【0183】

ブロック2018に戻って、検索値と対応する選択されたノード210の子孫がない場合、検索値は、インデックス200内で対応するエントリ144を持たず、インデックス不足としても言及される(ブロック2022)。該インデックス不足は、インデックストラバーサルを始めたキーマップ操作の種類に従い処理されてもよい(ブロック2024)。例えば、キーマップエントリ取得操作は、要求するクライアントに不足の適切な状態表示標示を返すことによって処理してもよい。対照的に、キーマップエントリ格納操作は、選択されたノード210の子孫としてインデックスに保存されるエントリ144と対応する新規ノード210を挿入することによって、インデックス不足を処理してもよい。例えば、新規ノード210が作成されてもよく、その様々なフィールドが、保存されるエントリ144のために適切に設定されてもよく、新規ノード210に対するポインタ218が、選択されたノード210内に保存されてもよい。新規ノード210がインデックス200に付加された場合、又は存在するノード210が修正された場合、付加された、又は修正されたノード210のすべての先祖ノード210のカウントフィールド214及びフィンガープリントフィールド216は、変更を反映するために更新されてもよい。

#### 【0184】

ブロック2006に戻って、選択されたノード210のタグ212が検索値と一致せず、検索値のプレフィックスである場合、インデックス不足が生じる場合があり、ブロック2022から処理が続けられる。一部の例において、この事例は、選択されたノード210がインデックス200のルートノードである時に生じる場合がある。同様に、一実施形態において、この新規ノード210をインデックス200に付加することは、検索値及び存在するルートノード210(この場合、選択されたノード210)のタグ212の両方に共通のプレフィックスであるタグ212を有する新規ルートノード210を作成することを含む。(一部の例において、新規ノード210の共通のプレフィックスは、あらゆる値のための有効なプレフィックスとして解釈される、空値である場合がある。)新規ルートノード210は、子孫として、選択されたノード210を参照するように設定されてもよい。必要な場合、付加的なノード210は、検索値と対応し、新規ルートノード210の付加的な子孫として構成されるために作成されてもよい。

#### 【0185】

一の実施形態において、インデックスの不足は、選択されたノード210のタグ21

10

20

30

40

50

2 が検索値と一致せず、検索値のプレフィックスでない場合、階層的な非平衡インデックス 200 をトラバースする間、ただちに生じなくてもよいことに留意されたい。一実施形態において、これらの事例に遭遇した場合、選択されたノード 210 が親を有する場合、親ノード 210 が選択される。親ノード 210 が別のインデックス 200 を参照する非内部ノードである場合、参照されたインデックス 200 のルートノード 210 が選択されてもよく、処理はブロック 2004 から続いてよい。そうでなければ、インデックスの不足が生じる場合がある。(しかしながら、この事例は他のインデックス 200 をインデックス化しない非階層的、内蔵型インデックス 200 には生じなくてもよいことに留意されたい。)この事例として、検索値が「alice」である図 19 の階層的インデックスを考慮されたい。インデックス 200a のトラバースが、タグ「ali」を有するノード 210w に処理されてもよい。ノード 210w は、「ali」と共に検索値のプレフィックスを形成する、関連するタグ情報「c」をもつ子孫ノード 210x に対するポインタを有するため、ノード 210x が選択されてもよい。しかしながら、ノード 210x のタグは、一致しない、検索値のプレフィックスでない、「alicia」である。したがって、トラバースは、インデックス 200c を参照する非内部ノードである、ノード 210w (ノード 210x の親)に戻ってもよい。同様に、トラバースは、ノード 210k に対して、そして最終的に検索値と一致するタグ 212 を有する、ノード 210m に対して続いてよい。

#### 【0186】

様々な実施形態において、非平衡インデックス 200、又は階層的な非平衡インデックス 220 は、キーマップインスタンス 140 内のキーマップエントリ 144 をインデックス化するために使用されてもよい。例えば、階層的インデックス 220 は、1 つ以上のパーティションインデックス 410、ブロックインデックス 420、又はエントリインデックス 430、又はキーマップインスタンス 140 内に実装される、あらゆる他のインデックスのレベルを実装するために用いられてもよい。上記に論考のように、異なるキーマップインスタンス 140 は、緩やかな同期プロトコルが用いられる場合、普通の操作の方法において相違、又は不一致になってもよい。一部の実施形態において、キーマップインスタンス 140 は、インデックス構造、又はインデックス化された内容の相違を識別するために、一貫した順序で(例えば、縦型、又は横型検索順序)、それぞれのインデックスデータ構造の各ノードをトラバースする、包括的なプロトコルを使用して同期化されてもよい。しかしながら、インデックスデータ構造内の、キーの数、及びカウントの含有及び/又は累積したハッシュ情報よりは、むしろキー情報に従ったデータの分散などの、上記に記載される非平衡インデックスの様々な特徴は、さらにコンピュータ的に効率的な同期化アルゴリズムの実装を容易にしてもよい。

#### 【0187】

前に説明したアンチエントロピー設定の調整プロトコルの多くの可能なバージョンは、キーマップインスタンス 140 によって実装された、非平衡で、多分満たされたインデックスの使用を意図する。このようなプロトコルの一実施形態の説明は、例えば、他の事例にわたりある事例を最適にするために選択することで、又はプロトコルの一般的なステップを行うために、アルゴリズムの 1 つ、又は別の特定の種類又はクラスを使用することで、一般的なプロトコルの予期される変化は、異なる実装の優先順位を提示し得ることが既知であるが、以下のとおりである。したがって、説明された実施形態は、制限することよりはむしろ事例として意図されたものである。

#### 【0188】

一実施形態において、非平衡インデックス 200、又は階層的な非平衡インデックス 220 の異なるインスタンスを調整するように構成されるアンチエントロピープロトコルは、様々な種類のメッセージのインスタンス間の交換を含んでもよい。アンチエントロピープロトコルの一実施形態に基づいてもよい例示的なメッセージの 1 組は、DATA メッセージ、REQUEST メッセージ、HASH メッセージ、FILTER メッセージ、及び FINGERPRINT メッセージを含んでもよい。これらの各メッセージのそれぞれの実



施形態の一般的な機能は、メッセージが、アンチエントロピープロトコルの実施形態を実装するためにどのように使用されてよいかの論考に続いて、以下に説明される。以下の論考において、当該キーマップインスタンスが、上記に記載された特徴のどれかを含む、非平衡インデックス 200 の 1 つ以上のインスタンス、又は層別非平衡インデックス 220 を実装してもよいことが理解されるが、参照が、キーマップインスタンス 140 間のデータの交換になされてもよい。

【0189】

D A T A メッセージは、1 つのキーマップインスタンス 140 から別のキーマップインスタンスへ、1 つ以上のインデックスノード 210 においてのデータを伝えるために使用されてもよい。一実施形態において、D A T A メッセージは、所定のノード 210 と関連するタグ 212 のみを伝えるように設定されてもよいが、他の実施形態において、D A T A メッセージは、所定のノード 210 と関連する他のフィールドを伝えてもよい。一部の実施形態において、所定のノード 210 が非内部ノードである場合、D A T A メッセージは、所定のノード 210 と関連するデータアイテムのすべて、又は一部分を含んでもよい（例えば、エントリ 144、又は別のインデックス 200 のルートノード 210 についての情報）。

【0190】

H A S H メッセージは、所定のノード 210 のフィールド、又は所定のノード 210 と関連するデータアイテムを明示的に伝えることなく、1 つのキーマップインスタンス 140 から別のキーマップインスタンスへ、1 つ以上のインデックスノード 210 についての情報を伝えるために使用されてもよい。一実施形態において、H A S H メッセージは、適切なアルゴリズムに従い計算された所定のノード 210 のハッシュのみならず、所定のノード 210 と関連するタグ 212 を伝えるように設定されてもよい。一部の実施形態において、所定のノード 210 のハッシュは、所定のノード 210 と関連するデータアイテム（例えば、キーマップエントリ 144）を反映してもよいが、所定のノード 210 のあらゆる子孫を除外してもよい。

【0191】

R E Q U E S T メッセージは、1 つ以上のノード 210 と関連する情報のための要求を伝えるために使用されてもよい。一実施形態において、R E Q U E S T メッセージは、1 つ以上のタグプレフィックス値を伝えるように設定されてもよい。それに応じて、要求するインスタンスは、伝えられたタグプレフィックス値が実際にプレフィックスである、タグ 212 を有するこれらのノード 210 についての情報を受け取ることを期待してもよい。所定のノード 210 について、受信した情報は、所定のノード 210 の対応するフィールドの内容及び/又は所定のノード 210 と対応するデータアイテム（例えば、キーマップエントリ 144）を含んでもよい。一部の実施形態において、特定のタグプレフィックス値によって定義された結果空間内の値、又は値の範囲などは、タグプレフィックス値のために返された結果から除外されるべきであることを特定することによって、R E Q U E S T メッセージは、要求されたプレフィックス値のさらなる必要条件をサポートしてもよい。例えば、R E Q U E S T メッセージは、プレフィックス「a l e x e」又は「a l e x j」とこれらのノード 210 を除き、タグプレフィックス値「a l e x」と一致するすべてのノード 210 についての情報が返されるべきであることを特定してもよい。

【0192】

今説明したメッセージは、個々のノード 210 の粒度のレベルで一般的に行われてもよい。しかしながら、キーマップインスタンス 140 間の相違が概して小さい場合（例えば、少数のノード 210 に伝えられた）、同期化過程が、一度に複数のノード 210 の状態を素早く確かめるために容易にしてもよい。一実施形態において、F I N G E R P R I N T 及び F I L T E R メッセージは、ノード 210 の総計についての情報を伝達するように設定されてもよい。とりわけ、一実施形態において、F I N G E R P R I N T メッセージは、1 つのキーマップインスタンス 140 から別のキーマップインスタンスへ、タグ 212 とともに、ノード 210 のフィンガープリントフィールド 216 を伝えるように設定さ

10

20

30

40

50

れてもよい。上記に記載されるように、所定のノード210のフィンガープリントフィールド216は、所定のノード210の子孫の機能として決定されるハッシュ値を保存するように設定されてもよい。したがって、異なるキーマッピングインスタンス140のそれぞれのノード210のフィンガープリントフィールド216が等しい場合、それぞれのノード210の子孫の配列及び内容が同一であることがほぼ確実であり得る（使用されたハッシュアルゴリズムの特徴によって）。すなわち、それぞれのノード210から下に向かっていくキーマッピングインスタンス140の部分が同期化されることは、ほぼ確実である。

#### 【0193】

フィンガープリントの使用は、相当数のノード210を含むキーマッピングインスタンスの140が同期化されたか否かの素早い判断を可能にする。しかしながら、対応する一部が同期化されていないことを示しているフィンガープリントは、概してどのように該部分が異なるかについてのさらなる詳細を提供しなくてもよい。一実施形態において、FILTERメッセージは、第1のキーマッピングインスタンス140から、第2のキーマッピングインスタンス140へ、特定のプレフィックス値に対応するノード210の数をコード化するフィルタ値を伝えるように設定されてもよい。第2のインスタンスは、ある場合は、どの第2のインスタンスのノード210が、第1のインスタンスに存在していないかを確認するために、プレフィックス値に対応する自体のノード210をテストするための受信したフィルタ値を使用してもよい。

#### 【0194】

一実施形態において、データ値の1組をフィルタ値に回復可能にコード化するあらゆる適切なフィルタリング技術が採用されてもよいことが考慮されるが、FILTERメッセージから伝えられたフィルタ値は、Bloomフィルタであってもよい。一般的に言えば、値の1組のBloomフィルタ（例えば、ノード210）は、Mが整数である、Mビットと対応してもよい。Bloomフィルタに、あらゆる値がコード化される前は、最初の値はゼロであってもよい。すなわち、フィルタのすべてのビットは、アサート停止状態であってもよい。Bloomフィルタは、それぞれが範囲[0、M-1]内の値でコード化されるべき値をマップする、フィルタ内でコード化すべき各値を、k独立ハッシュ機能の各1組に通すことによって投入されてもよい。各kの結果として得られるハッシュ値について、Bloomフィルタ内で対応するビットはアサートされる（例えば、論理1値に設定する）。M及びkは、偽陽性（以下に論考）の所望の可能性のみならず、Bloomフィルタ内でコード化されるべき値の数及び種類に従い、設計パラメータとして選択されてもよい。例えば、8のハッシュ機能を使用する1024ビットBloomフィルタにおいて、各ハッシュ機能は、アサートされるべきフィルタの1024ビットの特定の1つを特定する対応する10ビットハッシュ値を生成してもよい。

#### 【0195】

所定の値がBloomフィルタにコード化されたかどうかをテストするために、値は、フィルタをコード化するために使用されたk独立ハッシュ機能の同一の組を通り、フィルタ値の結果として得られたkビットが調査される。フィルタの結果として得られたkビットのどれもアサートされない場合、テスト値はフィルタ内で全くコード化されない。すべてのフィルタの結果として得られたkビットがアサートされる場合、テスト値はフィルタ内でコード化されても、されなくてもよい。すなわち、テスト値は、フィルタ内で最初にコード化されているか、又は偽陽性であってもよい。一部の実施形態において、ハッシュ機能は、所定の値の1組が、フィルタ内で無事にコード化される時に同一の偽陽性値が作成される可能性を低くするために、Bloomフィルタが作成された各別々の場合において無作為に、又は自然発生的に作成された（例えば、現在のシステム時間の機能として）ソルト、又はシード値でパラメータ化されてもよい。

#### 【0196】

したがって、例えば、第1のキーマッピングインスタンス140は、Bloomフィルタの中のプレフィックスPと対応するノード{A, B, C, D, E}の組をコード化してもよく、FILTERメッセージを使用する第2のキーマッピングインスタンス140に対してフ

フィルタを伝えてもよい。第2のキーマップインスタンス140において、ノード{A, B, X, Y, Z}の1組は、プレフィックスPと対応してもよい。第2のキーマップインスタンス140は、フィルタに対する各ノードをテストしてもよく、ノードA、B、及びXがフィルタでコード化されたが、ノードY及びZはフィルタで全くコード化されなかったことを決定してもよい。したがって、第2のキーマップインスタンス140は、第1のキーマップインスタンス140にノードY及びZが存在しないことを正しく結論づけてもよく、ノードA、B及びXが、Xが偽陽性である、第1のキーマップインスタンス140におそらく存在することを結論付けてもよい。その結果、第2のキーマップインスタンス140は、ノードY及びZについての情報を、第1のキーマップインスタンス140に伝えるための行動をとってもよい。

10

#### 【0197】

DATA、HASH、REQUEST、FINGERPRINT及びFILTERメッセージが実装され、あらゆる適切なプロトコル、又はAPIに従い伝えられてもよく、メッセージを解読及び適切に処理するために必要なあらゆる付加的な情報のみならず、上記に記載されるように、情報を伝えるように構成されるフィールド、又はパラメータの様々な種類を含んでもよいことを考慮されたい。一実施形態において、メッセージは、メッセージに含まれる所定のタグ値について、送信するキーマップインスタンスが、それぞれ、得られたデータ、及び必要なデータパラメータとして言及される、対応するデータ、又は必要な対応するデータのいずれかを有するかどうかを示す、付加的なパラメータを含んでもよい。例えば、キーマップインスタンス140が、タグ「a1」と、いくつかの数の子孫を有するノード210のためにFINGERPRINTメッセージを送信する場合、インスタンスは、「a1」によって定義されたプレフィックス空間内にいくつかのノード210を有することを示す、得られたデータパラメータを含んでもよい。またインスタンスは、例えば、「a1」によって定義されたプレフィックス空間のコピーが、不完全であると考えられる場合、必要なデータパラメータを含んでもよい。一部の実施形態において、DATA又はHASHメッセージは、必要なデータパラメータを明示的に特定してもよい一方、FILTER又はREQUESTメッセージは明示的に得られデータパラメータを特定するが、得られたデータパラメータは、DATA及びHASHメッセージに潜在してもよいが、必要なデータパラメータは、FILTER及びREQUESTメッセージに潜在してもよい。一実施形態においてFILTERメッセージは、必要なデータパラメータ、又は得られたデータパラメータの少なくとも1つを特定することが要求されてもよい。

20

30

#### 【0198】

一実施形態において、2つのキーマップインスタンスによって実行されたアンチエントロピープロトコルは、2つのインスタンスが互いに接点を確認する時に開始してもよい。各インスタンスは、その両方がいくつかのデータを持ち、欠くことを仮定してもよい。同様に、各インスタンスは、インスタンスのルートノード210のタグ212及びフィンガープリント216を特定し、得られたデータ及び必要なデータパラメータを含む、FINGERPRINTメッセージを他のインスタンスへ送信してもよい。例えば、層的非平衡インデックス220を用いるキーマップインスタンス140の実施形態において、ルートノード210は、親ノード、又は上位のインデックス200を持たないインデックス200内に、親ノードを持たないノード210と対応してもよい。

40

#### 【0199】

FINGERPRINTメッセージ処理の方法の一実施形態を、図21に示した。図示した実施形態において、操作は、FINGERPRINTメッセージが、メッセージ送信側から受信されたブロック2100で開始する。例えば、第1のキーマップインスタンス140は、タグ値、フィンガープリント、及び1つ以上の得られたデータ、又は必要なデータパラメータを含むFINGERPRINTメッセージを、第2のキーマップインスタンス140に伝えてもよい。FINGERPRINTメッセージが受信された後、メッセージ受信部のインデックスは、受信したタグ値が、対応するタグフィールド212のプレフィックスである(又は正確に一致する)ノード210が存在するかどうかを識別するた

50

めに、トラバースされてもよい（ブロック 2102）。例えば、キーマップインスタンス 140 のインデックスは、図 20 の方法を使用する、又は本願の適切な改良型を使用するルートノードから始まってトラバースされてもよい。

#### 【0200】

受信したタグ値がプレフィックスでない、又はいかなるノード 210 のタグフィールド 212 の正確な一致でない場合、FINGERPRINT メッセージによって参照されたノードと対応するノード 210 は、メッセージ受信部に存在しなくてもよい。同様に、受信部は、最初に受信した FINGERPRINT メッセージに含まれるタグ値を特定する送信側に、REQUEST メッセージをメッセージ送信側に伝えることによって応答してもよい（ブロック 2104）。一実施形態において、REQUEST メッセージの処理は、以下に詳しく説明するように続行されてもよい。一部の実施形態において、REQUEST メッセージは、受信された FINGERPRINT メッセージが得られたデータパラメータを示す場合のみに伝えられてもよい。

#### 【0201】

一部の実施形態において、アンチエントロピープロトコルの操作の間に交換された個々のメッセージの完了は、付加的なメッセージが、所定のメッセージが無事に完了することに応じて作成されたかどうかにより左右される。すなわち、一部の実施形態において、個々のメッセージを処理することは、他のメッセージに関して、処理状況を把握しない、非同期方法で生じてよい。本願に記載される例示的な実施形態の論考において、この処理状況を把握しない、非同期モデルが仮定される。したがって、REQUEST メッセージが作成された後、FINGERPRINT メッセージの処理は、それ自体が完了と見なされてもよい（ブロック 2106）。しかしながら、このモデルは、アンチエントロピープロトコルの一般的な操作にとって絶対に必要ではなく、他の実施形態において、いかなる所定のメッセージは、下位に作成されたメッセージとの、又は所定のメッセージへの応答において、同期化を遮断、待機、又は維持してもよい。例えば、明示的なハンドシェイク、肯定応答、再試行、又は他の種類のプロトコルは、1つのメッセージから別のメッセージへの完了の状態を伝えるために、一部の実施形態に用いられてもよい。

#### 【0202】

受信されたタグ値が、メッセージ受信部で特定のノード 210 のタグ 212 のプレフィックス、又は一致として対応する場合は、受信されたフィンガープリント値は、2つのフィンガープリントが一致するかどうかを決定するために、特定のノード 210 のフィンガープリントフィールド 216 と比較されてもよい（ブロック 2108）。一致する場合、メッセージ送信側とメッセージ受信側が、受信されたタグ値に関して同期化されることがほぼ確実（例えば、異なるデータから作成されているにもかかわらず、衝突する、又は同一の値を有する、2つのフィンガープリントを作るために使用する、フィンガープリントアルゴリズムの可能性に従い）である。例えば、プレフィックスとしての受信されたタグ値を有するいかなるノード 210 が、FINGERPRINT メッセージが送信されたキーマップインスタンス 140 及び、メッセージが受信されたキーマップインスタンス 140 内と同一の状態である。したがって、FINGERPRINT メッセージに応じて作成された付加的なメッセージはなくてもよく、メッセージは完了と見なされる（ブロック 2106）。

#### 【0203】

フィンガープリントが一致しない場合、メッセージ送信側とメッセージ受信側は、受信されたタグ値に関して同期化しておらず、送信側と受信側の状態が共に近くなるようにするためには付加的な作業が必要とされる場合がある。上記に記載されるように、FILTER メッセージは、送信側があるノード 210 についての特定の情報を受信側と通信することを可能にするのに役立つ。しかしながら、一部の実施形態において、妥当な偽陽性割合を保ちながら FILTER メッセージにコード化されてもよいノード 210 の数は、あるしきい値に限られてもよい。子孫ノード 210 の数が、受信されたタグ値と一致するメッセージ受信側ノード 210 でしきい値を超える場合、FILTER メッセージ

を送信する前に付加的な FINGERPRINT メッセージの処理を行うことはさらに効率的であり得る。

【0204】

したがって、図示した実施形態において、フィンガープリントが一致しない場合、メッセージ受信側の特定のノード 210 のカウントフィールドが、FILTER メッセージ処理のしきい値を超えるかどうか決定するために調査される (ブロック 2110)。超える場合、メッセージ受信側は、受信されたタグ値がプレフィックスである特定のノード 210 の子に従い、受信されたタグ値と対応するインデックス範囲の一部を細分化するように設定されてもよい (ブロック 2112)。各子ノード 210 について、メッセージ受信側は、それぞれの子ノード 210 のタグ 212 及びフィンガープリントフィールド 216 を特定して、最初のメッセージ送信側に、対応する FINGERPRINT メッセージを送り返すように設定されてもよい (ブロック 2114)。さらに、例えば、特定のノード 210 の子によって示されたような、受信されたタグ値と対応するインデックス範囲の一部に格差がある場合、メッセージ受信側は、格差と対応するタグ値にたいして 1 つ以上の REQUEST メッセージを送るように設定される (ブロック 2116)。受信された FINGERPRINT メッセージの処理は、そこで完了と見なされてもよい (ブロック 2118)。一実施形態において、上記の行動に加え、受信されたタグプレフィックス値が特定のノード 210 の完全な一致である場合、特定のノード 210 に対応する HASH メッセージは、メッセージ送信側に返されてもよい。

【0205】

例えば、図 17 に示すように、メッセージ受信側のインデックス 200 の特定のノード 210 a は、タグ「a l」、及び対応するタグ「a l a n」、「a l e x」、「a l f r e d」、「a l i」及び「a l z」を有する子を有してもよい。これは、メッセージ受信側が、「a l b」、「a l c」又は「a l d」で始まるであろうノード 210 についてではなく、「a l a n」及び「a l e x」で始まるノード 210 についてのある情報を有することを示唆する。同時に、メッセージ受信側は、子のタグ間の格差に対する REQUEST メッセージのみならず、ノード 210 a の各ノードに対する FINGERPRINT メッセージを伝えてもよい。負の REQUEST 構文がサポートされる実施形態において、メッセージ受信側は、特定のノードの子と対応するタグ以外のタグに対する REQUEST メッセージを伝えてもよい。例えば、メッセージ受信側は、「a l a n」、「a l e x」、「a l f r e d」、「a l i」及び「a l z」とプレフィックスであるタグ以外のタグに対して REQUEST メッセージを送ってもよい。

【0206】

特定のノード 210 のカウント値が処理している FILTER メッセージのしきい値を超えない場合、また、受信された FINGERPRINT メッセージが、得られたデータパラメータを含む場合、メッセージ送信側は、メッセージ受信側に存在しないノード 210 についての特定の情報を有してもよい。同様に、メッセージ受信側は、特定のノード 210 の子孫である各ノード 210 をフィルタ (例えば、上記記載の Bloom フィルタ) にコード化する、FILTER メッセージを送るように設定されてもよい (ブロック 2120 ~ 2122)。例えば、図 17 に言及して、特定のノードが、ノード 210 1 と対応する場合、各ノード 210 m ~ q をコード化する Bloom フィルタは、FILTER メッセージを介して作成され、返されてもよい。図示した実施形態において、得られたデータパラメータが最初の FINGERPRINT メッセージに含まれていなかった場合、それぞれの FINGERPRINT メッセージは、FILTER メッセージの代わりに、特定のノード 210 の各子に対して作成され、メッセージ送信側に返されてもよい (ブロック 2124)。これらの FINGERPRINT メッセージは、得られたデータパラメータを含んでもよい。この事例において、以下の FILTER 又は FINGERPRINT メッセージの作成のいずれかは、受信された FINGERPRINT メッセージの処理が完了してもよい (ブロック 2118)。

【0207】

F I L T E Rメッセージ処理の方法の一実施形態を、図 2 2 に示す。図示された実施形態において、操作は、例えば、上記に記載されるように、F I N G E R P R I N Tメッセージに応じて、タグ値及びフィルタ値を含むF I L T E Rメッセージがメッセージ送信側から受信されるブロック 2 2 0 0 で始まる。F I L T E Rメッセージが受信されると、メッセージ受信側のインデックスは、図 2 1 に関して上記に記載と同様の方法で、受信されたタグ値（例えば、受信されたタグ値がプレフィックスである、又は一致する）と対応するノード特定のノード 2 1 0 を識別するためにトラバースされる（ブロック 2 2 0 2 ）。一部の実施形態において、F I L T E Rメッセージが別のメッセージに応じて作成される場合、受信されたタグ値と対応するノード 2 1 0 は概して存在する。

#### 【 0 2 0 8 】

メッセージ受信者は、それがあある場合には、フィルタ値にコード化されないノード 2 1 0 を識別するために、F I L T E Rメッセージに提供されたフィルタ値に対する特定のノード 2 1 0 の各子孫をテストしてもよい（ブロック 2 2 0 4 ）。フィルタ値にコード化されないメッセージ受信側の各ノード 2 1 0 について、対応するD A T Aメッセージが、メッセージ送信側に返されてもよい（ブロック 2 2 0 6 ）。F I N G E R P R I N Tメッセージの処理は、そこで完了と見なされてもよい（ブロック 2 2 0 8 ）。上記に記載されるように、F I L T E Rメッセージに用いられるフィルタアルゴリズムの種類及び形態によって、偽陽性が生じる場合がある。すなわち、メッセージ受信側は、実際にはそうでない時に、ノード 2 1 0 のあるものがフィルタ値にコード化されており、したがって、メッセージ送信側と同一の状態を呈すると誤って結論付ける場合がある。したがって、アンチエントロピープロトコルの 1 ラウンドは、各ノード 2 1 0 に関して同期化になる 2 つのキーマッピングインスタンス 1 4 0 の結果にならなくてもよい。しかしながら、多くの実施形態において、アンチエントロピープロトコルの 1 ラウンドは、インスタンスがさらに相違するようになることを引き起こさなくてもよく、異なるインスタンス、及び使用されるアルゴリズムの特徴の程度によって、（例えば、フィルタのコード化のためのしきい値に定められた偽陽性の確率）プロトコルの反復のアプリケーションが、可能性のある程度を伴ういくらかのラウンド内で収束することが期待されてもよい。

#### 【 0 2 0 9 】

一部の実施形態において、H A S H、R E Q U E S T、及びD A T Aメッセージの処理は、F I L T E R及びF I N G E R P R I N Tメッセージよりもかなり単純である。一実施形態において、H A S Hメッセージ受信側は、メッセージに含まれるタグ値と対応するノード 2 1 0 を識別する試みをしてよく、識別されたノード 2 1 0 の対応するハッシュ値を計算してもよい。受信されたハッシュ値が、計算されたハッシュ値と一致する場合、識別されたノード 2 1 0 は、メッセージ送信側で対応するノード 2 1 0 とすでに同期化していてもよい。そうでない場合、受信されたタグ値を含むR E Q U E S Tメッセージは、送信側がさらに最新のデータバージョンを取得するために返される。

#### 【 0 2 1 0 】

R E Q U E S Tメッセージの処理は、一実施形態において、例えば、上記に記載される非平衡インデックス誘導技術を使用して、メッセージに含まれた受信されたタグ値が一致する、又は対応するタグ値 2 1 2 のプレフィックスである各ノード 2 1 0 を識別するメッセージ受信側を、単に含んでもよい。各識別されたノード 2 1 0 について、上記に記載されるように構成された、対応するD A T Aメッセージは、メッセージ送信側に返されてもよい。一実施形態において、受信されたD A T Aメッセージの処理は、メッセージ受信側で、メッセージに示されたタグ値と対応するノード 2 1 0 が存在するかどうかを識別することを含んでもよい。存在しない場合、対応するノード 2 1 0 が、メッセージから抽出されたデータと共に作成され、投入されてもよい。存在する場合、存在するノード 2 1 0 及び/又はその対応するデータ値に関連するデータは、メッセージから抽出されたデータと置き換えられてもよい。一部の実施形態において、存在するノード 2 1 0 のデータは、受信されたデータがさらに最新のものである場合にのみ置き換えられてもよい。例えば、D A T Aメッセージは、メッセージ送信側でノード 2 1 0 と対応するエントリ 1 4 4 の内容

10

20

30

40

50

を含んでもよく、エントリ 1 4 4 は、受信されたエントリ 1 4 4 が存在するエントリ 1 4 4 よりもさらに最新であるかどうかを確かめるために、メッセージ受信側で対応するタイムスタンプ情報と比較してもよい、タイムスタンプ情報を含んでもよい。最新である場合、受信されたエントリ 1 4 4 は、存在するエントリ 1 4 4 を置き換えてもよい。

#### 【 0 2 1 1 】

図 2 1 の一般的な同期プロトコルのバリエーションは可能であり、考えられる。例えば、一定の長さを有するパケットを使用して行われるキーマップインスタンス間の伝達の実施形態において、回線容量使用は、特定のノード 2 1 0 と対応する 1 つの F I N G E R P R I N T メッセージというよりはむしろ、1 パケット内の複数のノード 2 1 0 のための複数の F I N G E R P R I N T メッセージを伝えることによって、改善されてもよい。当該のパケットを受け取るインスタンスは、送信側とさらにメッセージを交換する必要なく、送信者と一致しないインデックス 2 0 0 の特定の 1 つを速やかに見分けることが出来てもよい。例えば、第 1 の F I N G E R P R I N T メッセージが一致しない場合、受信側は、パケットの送信側に対する R E Q U E S T、F I L T E R、又は他のメッセージを発行する前に、パケット内で、他の F I N G E R P R I N T メッセージを考慮してもよい。そうすることで、受信者は、相違をデータの特定の部分に絞ることが出来、すでに同期化されているデータ構造の他の部分に関するメッセージを交換するための不必要なネットワークトラフィックを削減してもよい。

#### 【 0 2 1 2 】

概して、アンチエントロピープロトコル及び / 又は更新伝播プロトコルを使用するキーマップインスタンス調整の実行のための上記に記載されるいずれかの方法、又は技術は、インスタンス内のキーマップインスタンス 1 4 0、又は個々のホスト 4 0 0 のレベルで操作されるように構成されるキーマップコーディネータ過程によって、実装されてもよいことを意図する。非平衡データ構造のためのアンチエントロピープロトコルを実装するための、前述の方法及び技術の多くの変更は可能であり、企図され、上記の論考は、限定的というよりは実例であることが意図される。例えば、他と頻繁に伝達する一部のエンティティを介するプロトコルの一般的なクラス、ネットワーク全体を通じて情報を分散するために無作為に選択されたエンティティは、ゴシップ系プロトコルとして言及されてもよく、ゴシップ系のプロトコルの他の技術、又は態様は、キーマップインスタンス 1 4 0 間のアンチエントロピープロトコルの使用に用いられてもよい。様々な実施形態において、上記に記載される実例同期化メッセージ（又は、他の適切なメッセージ）は、異なる特徴を有する同期プロトコルをもたらし異なる様式に組み合わせられてもよい。

#### 【 0 2 1 3 】

さらに、図 1 7 ~ 2 2 に関し、上記に記載される層的インデックス化データ構造及び同期化技術がキーマップインスタンス 1 4 0 内で使用するための効果的なデータ構造の実装の内容で論考されてきたが、当該データ構造及び同期化技術は、速やかなアクセスのためにインデックス化されてもよい、大量のデータのあらゆるアプリケーションに用いられてもよいことを意図したものである。当該のアプリケーションは、図 2 のシステムのような、オブジェクトストレージシステムを必ずしも含む必要はないが、データインデックス化が適用できる、データ系システム、検索システム、又はあらゆる他のアプリケーションを含んでもよい。

#### 【 0 2 1 4 】

様々な実施形態において、本願に説明された事例の無作為な作成、又は選択のいかなる種類の実装は、乱数、又は自称発生のためのいかなる適切なアルゴリズム、又は技術を用いてもよいことを留意したい。多くの場合、無作為の方法を実装するコンピュータ技術は、純粋に無作為な結果を生み出す、むしろ擬似乱数の結果を生む場合がある。例えば、擬似乱数アルゴリズムは、確率的に無作為な結果を作成するように構成される決定論的な過程を特定してもよい。本願で使用されるように、「無作為」又は「実質的に無作為」データの作成は、純粋に無作為データ供給源のみならず、あらゆる適切な擬似乱数コンピュータ技術を含むことを意図するものである。

## 【 0 2 1 5 】

ストレージサービスコンポーネント検出及び管理

ストレージサービスシステムの大規模な、極めて分散された実装において、システム全体にわたって分散された図 2 に示す多くの様々なシステムコンポーネントがあってもよい。例えば、ビットストアノード 1 6 0、コーディネータ 1 2 0、及びキーマップインスタンス 1 4 0 の数百、又は数千の例がある場合がある。このような規模の分散型システムの状態を管理することは、実用的な課題をもたらす。例えば、特定のシステムコンポーネントの異なるインスタンスは、計画されたメンテナンス、コンポーネントに依存するコンピュータ供給源の障害、機能に反してコンポーネントを孤立させる伝達障害、又は他の理由のため、いかなる時に稼動しない場合がある。さらに、新しい、又は以前のコンポーネントの非稼動は、任意の、又は予想不可能な時にある事例において稼動に戻る場合がある。

10

## 【 0 2 1 6 】

一実施形態において、発見、障害、及び検出デーモン ( D F D D ) 1 1 0 は、ストレージサービスシステムの様々な関連するコンポーネントの状態をそれぞれ監視するように構成され、当該の状態に関して互いに伝達しあうように構成され、当該のクライアントが、キーマップ、又はビットストア操作などの、システム操作を実行するために使用されてもよい、利用可能なシステムコンポーネントを識別してもよいことを通じるインターフェースで D F D D クライアントアプリケーションを提供するように設定されてもよい。概して、D F D D 1 1 0 は、他のコンポーネントに代わって、ストレージサービスシステムコンポーネントの現在の状態の様にアクセス可能なビューを提供するように設定されてもよい。すなわち、他の、類似しないコンポーネントとの状態情報の直接の伝達のために構成された、複数の異なるインターフェースを有するストレージサービスシステムの様々なコンポーネントを構成するよりはむしろ、当該の情報を提供し、それに依存する各コンポーネントは、基本的な D F D D インターフェースを介して D F D D 1 1 0 のインスタンスと伝達するように設定されてもよい。一部の実施形態において、D F D D 1 1 0 は、オペレーティングシステムによって管理される環境内で行われるように構成される、デーモン過程として実装されてもよい。しかしながら、他の実施形態において、D F D D 1 1 0 は、オペレーティングシステム、又は他のコンポーネントに依存、又は従属する必要なく、本願に記載される機能性を実装するように構成される、独立した、又は自立したハードウェア、又はソフトウェア媒体として、実装されてもよい。

20

30

## 【 0 2 1 7 】

一般的に言えば、D F D D 1 1 0 のインスタンスによって発見、及び監視されるように構成される、ストレージサービスシステムコンポーネントの各インスタンスは、アプリケーションインスタンスといわれることもある。例えば、操作状態、又は所定のビットストアノード 1 6 0 のヘルス状態は、所定のビットストアノード 1 6 0 によって実行されるために構成される S N M コントローラ 1 6 1 のインスタンスによって示されてもよい。したがって、S N M コントローラ 1 6 1 は、ビットストアアプリケーションインスタンスと対応してもよい。同様に、キーマップインスタンス 1 4 0 の操作状態は、キーマップインスタンス内の 1 つ以上のホスト 4 0 0 で実行されるために構成されるキーマップマネージャのインスタンスによって示されてもよい。各キーマップマネージャインスタンスは、キーマップアプリケーションインスタンスと対応してもよい。他の種類のアプリケーションインスタンスも可能であり、考えられる。例えば、一実施形態において、1 つ以上のストレージサービスシステムコンポーネントが展開されたものを介する各コンピュータシステムは、プロセッサ、メモリ、ディスク、入力/出力 ( I / O )、又は他のシステム供給源の使用などの、システム特定の操作状態の詳細を検出し、報告するように構成される、ホスト監視アプリケーションインスタンスを含んでもよい。一部の実施形態において、D F D D 1 1 0 の各インスタンスは、それ自体アプリケーションインスタンスとして設定されてもよい。すなわち、D F D D インスタンスは、他のアプリケーションインスタンスの状態に加えて、それ自体、操作状態を監視するように設定されてもよい。

40

## 【 0 2 1 8 】

50



ストレージサービスシステム内で、アプリケーションインスタンスは、一般的にアプリケーション名で識別され、それぞれのアプリケーションインスタンス識別子（ID）によって固有に識別されてもよい。例えば、特定のアプリケーション名は、「キーマップ・マネージャー」、「ビットストア・マネージャー」、「ホスト・マネージャー」、又は別の適切な名前などの、アプリケーションインスタンスの包括的な種類を識別する文字列を含んでもよいが、アプリケーションインスタンスIDは、アプリケーション名前領域内で特定のインスタンスを固有に識別する文字列を含んでもよい。一部の実施形態において、アプリケーションインスタンスIDは、「キーマップ・マネージャー - 4 A B 8 D 9 4 5」などの、アプリケーション名を明示的に含んでもよい。アプリケーションインスタンスIDのための他の適切な形式も用いられてもよい。一実施形態において、DFDD110の所定のインスタンスは、それぞれの状況情報を有する、多くのアプリケーションインスタンス（例えば、名前及びインスタンスIDを介して）を関連させるように設定されてもよい。例えば、図23に示す実施形態において、DFDD110は、それぞれがアプリケーション名112及びインスタンスID113とインスタンス状況情報114を結びつける、多くのエントリ111を含む。一部の実施形態においてDFDD110は、状態情報114と、所定のアプリケーション名112及びインスタンスID113との異なる種類の関連づけを反映するための1つ以上のテーブルを用いてもよいが、他の実施形態において、DFDD110は、ツリー、上記に記載されるような非平衡インデックス、又は所定のアプリケーションと対応する状態情報間の関連付けを暗示するあらゆる他の適切な種類のデータ構造を用いてもよい。

#### 【0219】

一部の実施形態において、アプリケーションインスタンスIDは、粒度の任意のレベルのそれ自体の名前領域を含んでもよいことに留意されたい。例えば、一実施形態において、所定のキーマップアプリケーションインスタンスIDは、〈マップ名〉/〈インスタンス〉/〈エンドポイント〉形式の場合がある。〈マップ名〉という用語は、所定のキーマップ展開と概して対応してもよい、キーエントリ関連付けの特定のキーマップ辞書を識別してもよい。（キーマップアプリケーションインスタンスが、異なるキーマップ展開をDFDD110の同一のインスタンス内で管理することは可能である。）〈インスタンス〉という用語は、固有の文字列によって、キーマップインスタンス140内の特定のホスト400を識別してもよい。〈エンドポイント〉という用語は、識別されたホスト400（例えば、特異的な過程）で実行する多くの機能的に同一のキーマップアプリケーションの1つを識別してもよい。アプリケーションインスタンスID内の他の複雑な名前領域は、可能であり、考えられる。

#### 【0220】

DFDD110によってアプリケーションインスタンスと関連する状態情報は、様々な異なる種類の情報を含んでもよい。一実施形態において、DFDD110は、DFDD110によって管理されるすべての種類のアプリケーションインスタンスに共通であってもよい、広域状態情報である、状態情報114内に保存されるように設定されてもよい。例えば、さらに以下に詳しく説明するように、一部の実施形態において、DFDD110は、状態の1組の間の考えられる移行のみならず、アプリケーションインスタンスの広域操作状況（又は、単に広域状態）の1組を定義する、広域操作状態マシンを実装してもよい。このような実施形態において、DFDD110によって管理された各アプリケーションインスタンスは、いかなる時に、広域状態の組の特定の1つと関連してもよく、所定のアプリケーションインスタンスのための広域状態は、アプリケーションインスタンスのマシン、及び動作に従い、時間と共に変化してもよい。

#### 【0221】

アプリケーションインスタンスの広く異なる種類に共通してもよい、広域状態情報に加えて、一部の実施形態において、状態情報114は、特定のアプリケーションインスタンス又はインスタンスの種類に特定されてもよく、又はカスタマイズされる場合がある操作状態情報を反映してもよい。例えば、アプリケーションインスタンスが、特定のビットス

トアノード１６０のビットストアマネージャーと対応する場合、その状態情報１１４は、その特定のノードで利用可能なストレージリソースの量、これらの供給源の該種類（例えば、高性能、低性能など）、又はビットストアノードの該内容に特有である、あらゆる他の関連する状態情報を含んでもよい。同様に、特定のキーマップインスタンス１４０のキーマップマネージャーに対応するアプリケーションインスタンスについては、その状態情報は、特定のキーマップインスタンス、使用される、又は利用可能なキーマップストレージリソース、又は、他の関連するキーマップ状態情報によって管理されるエントリ１４４の数についての情報を含んでもよい。一部の実施形態において、どのアプリケーションインスタンス特有の状態情報を対応するＤＦＤＤエントリ１１１内に含むかという選択は、ＤＦＤＤクライアントの要求に従い決定されてもよい。例えば、いくつかの選択肢から特定のビットストア又はキーマップアプリケーションを選択することにおいて、コーディネータ１２０又はノードピッカー１３０を援助するために使用できてもよい状態情報は、これらのアプリケーションインスタンスのＤＦＤＤエントリ１１１内に含まれてもよい。

10

#### 【０２２２】

一部の実施形態において、アプリケーションインスタンスの状態情報１１４は、どのようにＤＦＤＤクライアントがインスタンスへアクセスするかについての情報も含んでもよい。例えば、状態情報１１４は、ＤＦＤＤクライアントが、アプリケーションインスタンスとの連絡を通じて確立してもよい、インターネットプロトコル（ＩＰ）アドレス及びポート番号を含んでもよい。一部のアプリケーションインスタンスは、ウェブサービスインターフェース、出版／購読型インターフェース、又は他の適切なインターフェースなどの他の種類のインターフェースをサポートしてもよい。このような実施形態において、出版／購読チャンネルを購読するための、又はアプリケーションインスタンスとの通信を確立するために必要な別の種類の処置を行うために、状態情報１１４はＵＲＬ、又はＤＦＤＤクライアントがウェブサービスコールを実行するために必要な他の情報を含んでもよい。一部の実施形態において、アプリケーションインスタンスアクセス情報に加えて、又はその代わりに、状態情報１１４は、ストレージサービスシステム内に物理的に位置付けられている場所についての情報を含んでもよい。例えば、状態情報１１４は、特定のアプリケーションインスタンスが対応する、データセンタ３００、又は領域３１０の識別子を含んでもよい。

20

#### 【０２２３】

上記に記載されるように、一部の実施形態において、ＤＦＤＤ１１０は、所定のアプリケーションインスタンスが通常に操作しているか、したがって、使用できるか、又は異常な状態にあるかどうかを一般用語で示してもよい、個々のアプリケーションインスタンスのための広域状態情報を維持してもよい。一実施形態において、ＤＦＤＤ１１０のインスタンスによって監視するように構成された各アプリケーションインスタンスは、（必ずしもではないが）多くの場合、数秒、又は数分などの通常の間隔で、ＤＦＤＤ１１０へその状態を報告するように設定されてもよい。このような報告は「ハートビート」と言われることもある。ハートビート報告は、あらゆる適切なプロトコルに従い（例えば、ＴＣＰ／ＩＰメッセージとして、ウェブサービスコールとして、又は他の基準、又は所有者メッセージプロトコルに従って）通信されてもよく、情報の内容によって様々であってもよい。最小限の例として、所定のアプリケーションインスタンスは、所定のインスタンスに対応するアプリケーション名及びアプリケーションインスタンスＩＤを単に含むＤＦＤＤ１１０に対するハートビートを提出してもよい。他の事例において、所定のアプリケーションインスタンスは、ローカル供給源利用の特定の状態などの、ハートビートにおける付加的な状態情報を含んでもよい。一部の実施形態において、アプリケーションインスタンスは、ハートビートを送信する前に、それ自体の機能状態を確認するために、あるレベルの自己診断、又は自己照合を実行するように設定されてもよいが、他の実施形態において、アプリケーションインスタンスは、いかなる自己評価に依存することなく、ハートビートを送信してもよい。

30

40

#### 【０２２４】

50

一般的に言えば、予想通りに、アプリケーションインスタンスが、DFDD110へハートビートを送信している場合、それは通常に操作しているという理にかなった可能性がある。ハートビートがある期間にわたって中断されなければならない場合、アプリケーションインスタンスに異常があるという理にかなった可能性がある。図24は、ハートビートアクティビティ及び/又は他のパラメータの機能として、各アプリケーションインスタンスのためにDFDD110によって維持されてもよい広域状態マシンの一実施形態を示す。図示した実施形態において、新規のアプリケーションインスタンスは、NEW状態でオンライン化され、例えば、間もなく操作を開始し、DFDD110のインスタンスの存在を知らせ、アプリケーション名、アプリケーションインスタンスID、及び対応するエントリ111を使用するために、DFDD110に必要なあらゆる他の情報を提供する。新規のアプリケーションインスタンスが安定し、通常の操作を開始すると、OK状態に入る。様々な実施形態において、NEWからOK状態への移行は、時間の機能（例えば、アプリケーションインスタンスの種類に基づいた初期設定時間）、アプリケーションインスタンス自己報告、管理者介入、又はこれら、又は他の要素の組み合わせであってもよい。

#### 【0225】

図示した実施形態において、インスタンスのDFDD110への最後のハートビートが障害しきい値 $T_{fail}$ 未満であるため、アプリケーションインスタンスは、時間が経過した間、OK状態を維持してもよい。例えば、DFDD110は、対応するインスタンスから受信した各ハートビートに増加された各アプリケーションインスタンスのためのカウンタを維持してもよく、 $T_{fail}$ が経過する前にその値が変更するかどうかを確かめるためにそれぞれのカウンタ（例えば、カウントダウン時間）を監視してもよい。一部の実施形態において、以下に説明するように、当該の状態のなかで相違がある場合があるが、OK（及び、あるいはNEW）以外の広域状態は、異常な操作状態、又は障害状態として概して言及される場合がある。

#### 【0226】

アプリケーションインスタンスのための最後のハートビートから時間 $T_{fail}$ が経過した場合、その広域状態は、INCOMMUNICADOへ移行してもよい。図示した実施形態において、INCOMMUNICADOは、アプリケーションインスタンスに異常があるが、永久に障害があることを断定的に決定されていないことを示す過渡状態として機能してもよい。例えば、アプリケーションインスタンスは、一時的に行き詰まる、又はハングアップする場合があります、DFDD110へのハートビートメッセージは、遅延される、又は失われる場合があります、又は以下にさらに詳しく説明するように、DFDD110の一例は、アプリケーションインスタンスの現在の状態に関してDFDD110の別のインスタンスと同期化しない場合がある。ハートビートがINCOMMUNICADO状態のアプリケーションインスタンスから受信される場合、インスタンスはOK状態へ逆移行してもよい。一部の実施形態において、DFDDクライアントは、自身のリスクでINCOMMUNICADO状態にあるアプリケーションインスタンスを使用することを選んでよい。

#### 【0227】

アプリケーションインスタンスが、INCOMMUNICADO状態から自然に回復しない場合、インスタンスに影響しているさらに深刻な問題がある場合がある。図示した実施形態において、2つの予想される障害のシナリオが生じる場合がある。FAIL状態によって示すように、個々のアプリケーションインスタンスは、例えば、個々のインスタンスをホストしている根本を成すコンピュータ供給源の障害によって、分離に失敗する場合がある。あるいは、アプリケーションインスタンスは、NETWORK SPLIT状態によって示すように、インスタンスとDFDD110との間のネットワーク通信の喪失のために機能しない場合がある。例えば、アプリケーションインスタンスは、互いにストレージサービスシステムの部分を隔絶する通信障害のため、他ではなく、DFDD110の一部のインスタンスに対して操作可能であり、アクセス可能であってもよい。

#### 【0228】

所定のアプリケーションインスタンス障害が隔絶されている、又はネットワークスプリットのためあるかどうかを確実に決定するのは困難である場合がある。一部の実施形態において、DFDD110は、アプリケーションインスタンスが、INCOMMUNICADO状態からFAIL状態、又はNETWORK SPLIT状態へ移行すべきかどうかの決定をするための利用可能な情報の様々な種類を考慮に入れる、それぞれ発見的な基準 $H_{fail}$ 及び $H_{netsplit}$ を用いてもよい。例えば、該基準は、所定のアプリケーションインスタンスが、別の障害状態に移行する前に、少なくともしきい値量の時間 $T_{heuristic}$ にわたってINCOMMUNICADO状態であることを必要としてもよい。さらに、該基準は、所定のアプリケーションインスタンスとしての、同一の領域310、又はデータセンタ300と、供給源を共有する、又はそれに属する他のアプリケーションインスタンスもINCOMMUNICADO、FAIL又はNETWORK SPLIT状態であるかどうかを考慮に入れてもよい。例えば、所定のアプリケーションインスタンスとしての同一のIPアドレス上に、又は同一の領域310、又はデータセンタ300内の別のアドレス上に位置する別のアプリケーションインスタンスがOKである場合、所定のアプリケーションインスタンスの障害は隔絶されている可能性がある場合がある。それに反して、複数のアプリケーションインスタンスがOKではない場合、とりわけ、アプリケーションインスタンス状態が地理、又はネットワークトポロジーに従いクラスタ化された場合に、ネットワークスプリットのシナリオである可能性がある場合がある。一部の実施形態において、DFDD110は、障害の本質を決定するための消極的に受信した状態情報を使用することに加えて、障害の疑いのあるアプリケーションインスタンスを問い合わせるように設定されてもよい。一部の実施形態において、発見的な基準は、アプリケーションインスタンスが、あるしきい値の可能性（例えば、50%可能性より高い、90%可能性より高い、など）に従い、確率的に機能不全の可能性があるかどうかを決定するように設定されてもよい。

#### 【0229】

発見的な基準によって、機能不全のアプリケーションインスタンスは、FAIL状態、又はNETWORK SPLIT状態のいずれかに移行してもよい。一部の実施形態において、ハートビートが受信された場合、これらの状態のいずれかからOK状態に戻る移行をしてもよいが、他の実施形態において、これらの状態のいずれか、又はその両方は回復可能であってもよい。INCOMMUNICADO状態にあるアプリケーションインスタンスは、障害の可能性に機能的、又は回復可能であることが想定される場合があるが、FAIL又はNETWORK SPLIT状態にあるアプリケーションインスタンスは、機能不全であることが予想される場合がある（一部の実施形態において、回復の可能性を有する）。概して、DFDDクライアントは、これらの障害状態のいずれにおいても、これらのアプリケーションインスタンスを選択することを避けてもよい。一部の実施形態において、DFDD110は、これらの障害状態のいずれにおいてものアプリケーションインスタンスについての情報を、クライアントから隠すように設定されてもよい。

#### 【0230】

図示した実施形態において、アプリケーションインスタンスは、FORGOTTEN状態に進める前に、それぞれ $T_{clean}$ 及び $T_{recover}$ 期間に、FAIL又はNETWORK SPLIT状態にとどまってもよい。例えば、FAILのある事例において、機能不全のアプリケーションインスタンスと関連する供給源は、回復、又は分析の目的のために一定期間保存されてもよい。可能な場合、当該の供給源（例えば、ビットストアノード160のストレージリソース）は、新規のアプリケーションインスタンスとして再展開するために初期化されてもよい。NETWORK SPLIT状態のある事例において、機能不全のアプリケーションインスタンスのないシステム操作を続けるかどうかに関連する決定をする必要がある場合があり、機能不全がある場合、どのような回復処置がとられるべきかを決定する必要がある場合がある（例えば、残存するアプリケーションインスタンスの間でオブジェクトの複製を再生成することなど）。一部の実施形態において、機能不全のアプリケーションインスタンスは、当該の回復処置が完了するまで、FORGOTTEN状

態を通過しなくてもよい。

#### 【0231】

アプリケーションインスタンスのFORGOTTEN状態は、DFDD110内で明確に表されなくてもよい。むしろ、一部の実施形態において、DFDD110からのDFDDエントリ111などの、アプリケーションインスタンスの存在する状態情報の削除によってマークされてもよい。一般に、ある例において、アプリケーションの新規のインスタンスは、NEW状態を介してフォーガットンインスタンスへ割り当てられた同一の供給源を使用して初期化されてもよいが、アプリケーションインスタンスは、FORGOTTEN状態から回復しなくてもよい。一部の実施形態において、FORGOTTEN状態にある間に、アプリケーションインスタンスが自然にハートビートの送信を再開しなければならない場合、DFDD110は、インスタンスが忘れられた（例えば、有効なエントリ111ともはや対応しない）ことを認識してもよく、インスタンスに操作を中止する、又はリセットする、又はそれ自体を再初期化することを指示してもよい。

10

#### 【0232】

一部の実施形態において、広域状態の移行を考慮に入れる経験則及び移行時間のパラメータは、アプリケーションインスタンスの異なる種類に対して異なってもよく、これらのパラメータの一部、あるいはすべては、DFDDクライアントによって調整されてもよいことに留意されたい。また、DFDDクライアントは、概してDFDD110のインスタンスを所定のアプリケーションインスタンスの現在の広域状態を確かめるために問い合わせるが、一部の実施形態において、DFDD110は、出版／購読状態の変更通知モデルをサポートしてもよい。例えば、DFDDクライアントは、クライアントが、特定のアプリケーションインスタンス、又はインスタンスの組の広域状態変更のすべての、又はある種の通知を受けたいということを、購読過程を介してDFDD110へ通知してもよい。このような状態変更を検出する際に、DFDD110は、購読するDFDDクライアントに対して変更のメッセージ表示を伝えてもよい。

20

#### 【0233】

しばしば、アプリケーションインスタンスは、アプリケーションインスタンスに最も近いDFDD110のインスタンスへ、ハートビート情報を送信するように設定されてもよい。例えば、一部の実施形態において、アプリケーションインスタンスが、単にホストのローカルIPアドレスを参照すること及びアプリケーションインスタンスDFDD通信のために用意される既知のIPポートを使用することによって、DFDD110のローカルインスタンスへ直ちにアクセスできるように、DFDD110のインスタンスは、1つ以上の他のアプリケーションインスタンスをホストするように構成される各コンピュータシステムで提供されることがある。しかしながら、アプリケーションインスタンスが、他ではなく、DFDD110の一部のインスタンスへその状態を報告し、その状態を同期化するためのある取り組みがない場合、DFDD110の展開されたインスタンスは不一致になってもよい。

30

#### 【0234】

一部の実施形態において、DFDD110のインスタンス間の不一致は、ゴシップ系のプロトコルなどのキーマップインスタンス140に関する上記に記載されるものと類似の同期プロトコルを使用して対処されてもよい。しかしながら、多くの事例において、DFDD110のインスタンスによって集合的に管理されるDFDDエントリ111の数は、キーマップインスタンス140によって管理されるキーマップエントリ144の数よりも実質的に少なくてもよい。この事例の場合、簡素化された調整プロトコルがDFDD110のインスタンスを同期化するために使用されてもよい。このようなゴシップ系のプロトコルの一実施形態の操作の方法は、図25に示す。図示した実施形態において、開始しているキーマップインスタンスとしても言及される、DFDD110の1つのインスタンスであるブロック2500で始まる操作は、同期化のためにDFDD110の別のピアインスタンスをランダムに選択する。一部の実施形態において、開始するDFDDインスタンスは、開始しているDFDDインスタンスの状態情報に従い、現在機能不全状態にある（

40

50

例えば、NETWORK SPLIT)これらのDFDDインスタンスの中からピアDFDDインスタンスを時には故意に選択してもよい。開始しているDFDDインスタンスが明らかに機能不全であるピアインスタンスとの連絡及び同期化に成功する場合、明らかな機能不全からの回復が容易になってもよい。

#### 【0235】

開始しているインスタンスは、そのエントリ111に反映されるアプリケーションインスタンスの識別する情報のハッシュ値を計算してもよい(例えば、アプリケーションインスタンス名及びID、及びおそらくエンドポイント、又は他の識別する情報をハッシュすることによって)(ブロック2502)。ハッシュ値は、例えば、MD5アルゴリズムなどのあらゆる適切なハッシュアルゴリズムに従い決定されてもよい。開始しているインスタンスは、現在のアプリケーションインスタンス状態情報の保存されたリストと共に、ピアインスタンスへ、計算されたハッシュ値を伝えてもよい(例えば、ハートビートカウント、広域状態情報及び/又は状態情報114に含まれるあらゆる他の情報)(ブロック2504)。状態情報のリストは、開始しているインスタンス及びピアインスタンスの両方に一貫したリストを作成する、あらゆる基準に従い保存されてもよい。例えば、該リストは、アプリケーションインスタンス名及び/又はIDに従い保存されてもよい。

#### 【0236】

上記に記載されるように、一部の実施形態において、アプリケーションインスタンスと関連する状態情報は、ハートビートメッセージ内に含まれるハートビートカウントから派生してもよいことに留意されたい。同様に、一部の実施形態において、DFDDインスタンスは、アプリケーションインスタンスのためのハートビートカウント情報を交換してもよく、他のDFDDインスタンスから直接に状態情報を受け取るよりはむしろ、受信したハートビートカウント情報からアプリケーションインスタンスの状態情報を派生してもよい。したがって、一実施形態において、所定のDFDDインスタンスは、その情報が特定のアプリケーションインスタンスから直接に受信されたか、又は同期プロトコルを介して別のDFDDインスタンスから非直接的に受信されたかに関わらず、受信されたハートビートカウント情報に基づいて、特定のアプリケーションインスタンスの状態を(例えば、図24の状態マシンに従い)更新するように設定されてもよい。このような実施形態において、DFDDインスタンス間の状態情報を操作するアプリケーションインスタンスの同期化は、同期プロトコルの操作を簡略化してもよいアプリケーションインスタンスの特定の広域操作状態(例えば、OK、INCOMMUNICADOなど)を、直接交換することなく、ハートビート情報の同期化を伴ってもよい。

#### 【0237】

状態情報のハッシュ値及びリストの受信に応じて、ピアインスタンスは、開始しているインスタンスによって行われるものと一致する方法で、それ自体のアプリケーションインスタンスの識別する情報のハッシュ値を計算し(ブロック2506)、結果として得られたハッシュ値と、開始しているインスタンスから受信されたハッシュ値とを比較する(ブロック2508)。2つの値が一致する場合、開始しているインスタンスとピアインスタンスの両方が、アプリケーションインスタンスの同一の組と一致するエントリ111を有することが高確率である。ピアインスタンスは、状態情報の受信したリストを走査し、適切であれば受信したリストからそのエントリ111を更新する(ブロック2510)。例えば、受信したリストのハートビートカウント、又はタイムスタンプが、ピアエントリ111の1つに保存されたものよりも大きい、又はより最近である場合、該ピアは、受信されたリストの状態情報からエントリ111を更新してもよい。一部の実施形態において、ピアインスタンスは、開始しているインスタンスからのリスト受信と同時に、それに続いてのいずれかで、同様の処理のための開始しているインスタンスへ、状態情報の独自のリストを送り返してもよい。

#### 【0238】

ハッシュ値が一致しない場合、少なくとも1つのエントリ111で、ピア及び開始しているインスタンスに既知のアプリケーションインスタンスの組が異なることがあり得る。

同様に、ピアインスタンスは、開始しているインスタンスに既知のエントリ 1 1 1 の完全なダンプを要求してもよい（これらエントリ 1 1 1 の状態情報 1 1 4 にちょうど対立する形で）（ブロック 2 5 1 2）。ピアインスタンスは、欠けていたあらゆるエントリ 1 1 1 を加え、残存するエントリ 1 1 1 の状態を同期化してもよい（ブロック 2 5 1 4）。上記のように、一部の実施形態において、ピアインスタンスは、開始しているインスタンスからのダンプ受信と同時に、それに続いてのいずれかで、開始しているインスタンスへそのエントリ 1 1 1 の完全なダンプを送り返してもよい。

#### 【 0 2 3 9 】

一部の実施形態において、システム内に存在する D F D D 1 1 0 の各インスタンスは、一定間隔で、説明した同期プロトコル、又は本願の適切な改良型を繰り返し実行するように設定されてもよいことに考慮されたい。例えば、プロトコルは、およそ定期的に 1 秒間隔、又はあらゆる他の適切な間隔で、D F D D 1 1 0 のインスタンスによって実行されてもよい。さらに、一部の実施形態において、D F D D 1 1 0 のインスタンスは、およそ同一の間隔で同期プロトコルを実行してもよいが、その時々で、D F D D 1 1 0 のインスタンスの部分のみがプロトコルを開始する場合がある、異なる相は、互いに関連してオフセットであってもよい。

#### 【 0 2 4 0 】

一部の実施形態において、D F D D 1 1 0 のインスタンスは、単にストレージサービスシステム内に定義されているアプリケーションインスタンスではなく、あらゆる分散型システム内のあらゆる種類のアプリケーションインスタンスの状態情報を調整、及び通信する操作に使用される場合があることに留意されたい。また、一部の実施形態において、異なるグループの D F D D インスタンスは、異なるアプリケーションインスタンスの状態情報を管理してもよい。このような一部の実施形態において、グループは、同一グループのメンバーであり、D F D D 同期化の条件として、識別子が一致することを要求する D F D D 1 1 0 のインスタンスに共通識別子を割り当てることによって、互いを識別してもよい。例えば、ストレージサービスシステムのアプリケーションインスタンスを管理する D F D D インスタンスは、ストレージサービスシステムと関連していない、その他のアプリケーションインスタンスの状態を管理するように設定されている D F D D インスタンスと異なる識別子を有してもよく、同一識別子を有するそれらの D F D D インスタンスのみが、図 2 5 の同期プロトコルに従い、互いに情報を交換してもよい。

#### 【 0 2 4 1 】

一部の実施形態において、D F D D グループ識別子は、同一システム内に存在するアプリケーションインスタンスの異なる設定を識別するために使用されてもよい。例えば、「製品」識別子に対応する D F D D 1 1 0 のインスタンスのある組は、ストレージサービスシステム、又は別の分散型システムの製品バージョンを管理するために配置されてもよく、製品システムに対応するアプリケーションインスタンスの組を示す場合があり、一方、別の「テスト」識別子に対応する D F D D 1 1 0 のインスタンスの組は、アプリケーションインスタンス及び状態の異なる組に対応するシステムのテストバージョンを管理するために配置されてもよい。一部の例では、いずれかのシステムバージョンに対応するアプリケーションインスタンス及び/又は D F D D インスタンスは、同一の基礎システムリソース上（例えば、同一コンピュータシステム上）で実行されてもよいが、それらの異なる D F D D グループ識別子の長所によって、互いに透過的に行われる場合があることに留意されたい。例えば、図 2 5 に示されるプロトコルのような同期プロトコルの実行中に、D F D D インスタンスは、同一グループのメンバーであるかどうかを最初に判断し（例えば、グループ識別子を交換することによって）、この判断を条件とする後続の同期化ステップを実行し、従ってグループ間のアプリケーションインスタンスの状態情報の分離を促進してもよい。

#### 【 0 2 4 2 】

前述の D F D D 1 1 0 のインスタンスを同期化するためのゴシップ系プロトコルは、ストレージサービスシステム全域に渡る既存のアプリケーションインスタンスの動作状態の

配布を助長するだけでなく、その他のシステムコンポーネントによる新しいアプリケーションインスタンスの発見も促進する場合があることに留意されたい。例えば、一度新しいアプリケーションインスタンスが初期化され、DFDD110のインスタンスと連絡を取る（例えば、システム上で局所的に動作しているインスタンス上で新しいアプリケーションインスタンスが初期化された）と、新しいインスタンスに対応する新しいエントリ111が作成されてもよい。新しいエントリ111が作成されたDFDD110のインスタンスが、DFDD110の様々なその他のインスタンスでその状態を同期化する際、新しいエントリ111は、システム全体に伝播されてもよい。様々な目的（例えば、新しいオブジェクト30を保存する、又はキーマップエントリ140を更新する）のためのアプリケーションインスタンスを識別するために、DFDD110に問い合わせるDFDDクライアントは、あらゆる既存のものばかりでなく、新しいアプリケーションインスタンスに関する状態情報を提供されてもよい。

10

#### 【0243】

上記に記載される実施例において、障害検出及び発見に関するアプリケーションインスタンスの状態変更は、アプリケーションインスタンスの一部又はこれらのインスタンスを参照するDFDDクライアントの干渉なしに、システム全体に伝播されてもよい。つまり、所定のアプリケーションインスタンスは、どのようにハートビート情報をDFDD110の1つのインスタンスに伝達するかのみを知っていればよい。システム内のDFDD110のすべてのインスタンス、その他のアプリケーションインスタンス、又は当該の所定のアプリケーションインスタンスを起動する様々なクライアントの知識を有する必要はない。同様に、DFDDクライアントは、システム内のその他のクライアント、又はすべてのアプリケーションあるいはDFDDインスタンスの独立した知識を有する必要はなく、クライアントは、システム内で使用可能なリソースの状態に関する合理的に最新の情報を獲得するために情報をやり取りする、DFDD110のインスタンスに依存してもよい。その他のアプリケーションインスタンスを要求せずに、アプリケーションインスタンスの状態を変更すること、又はクライアントにそのような変更を即座に知らせることを可能にすることより、DFDD110は、ストレージサービスシステムの拡張性を促進してもよい。

20

#### 【0244】

#### ストレージクラス

30

ストレージサービスシステムの一部の実施形態において、オブジェクト30は、それらの複製レベル、領域310に渡る複製の分散、複製が保存されているストレージリソースの種類、及び/又はその他のシステム機能あるいはポリシーに関して一様に取り扱われる場合がある。例えば、システムは、同一回数だけ、同一番号の異なる領域310に、各オブジェクト30を複製することを試みる場合がある。しかし、異なるクライアント50は、異なるオブジェクト30に対して、異なるストレージ要件を有してもよい。例えば、1つのクライアント50は、初期ストレージポリシーが提供し得る信頼性の程度よりも高い信頼性で特定のオブジェクト30を保存する（例えば、数及び複製の分散に関して）ことを望む一方、別のクライアント50は初期設定レベルの信頼性でさえ必要としない場合がある。あるいは、クライアント50は、可能な信頼性の費用で、オブジェクトの複製が分散する領域310の数を制限することによって、オブジェクトの書き込み性能を向上することを希望する場合がある。

40

#### 【0245】

同様に、一実施形態において、図2のようなストレージサービスシステムは、オブジェクト30のストレージクラスをサポートするように設定されてもよい。一般的にいうと、所定のオブジェクト30のストレージクラスは、当該所定のオブジェクト30に関するサービスレベル合意書（SLA）に影響を与えるあらゆるストレージサービスシステム機能又は特徴の組を指定してもよい。サービスレベル合意書は、一般的に、クライアントから受け取る、ある対価（例えば、料金又はその他の対価の適切な種類）と引き換えに、サービスプロバイダがクライアントに提供するサービスに対する保証又は期待を反映してもよ

50



い。例えば、ストレージサービスシステムにより管理されるオブジェクト30のSLAは、様々なレベルのオブジェクトの信頼性、可用性、アクセス機能（例えば、待ち時間、バンド幅）、料金又はサービス率、あるいはクライアントのオブジェクト30とのやり取りにおける、あらゆるその他の測定可能な態様を指定してもよい。一部の実施形態において、ストレージクラスは、SLA特性（例えば、以下に記載されるオブジェクトの複製の数及び配布）の特定のサブセットのみを指定してもよいが、他の実施形態において、ストレージクラスは所定のオブジェクト30に関して、SLA合意のすべての定義された態様を網羅する総合的なSLAと直接対応してもよい。

#### 【0246】

一実施形態において、ストレージサービスシステムは、それぞれが特定の定義されたSLA特性を有する、ストレージクラスの固定組を定義する場合があり、クライアント50は、特定のオブジェクト30と特定のストレージクラスを関連付けることを選択してもよい。例えば、初期ストレージクラスは、オブジェクト30が少なくとも3回、少なくとも2つの異なる領域310に複製されるように指定してもよい。高信頼性ストレージクラスは、オブジェクト30が少なくとも5回、少なくとも3つの異なる領域310に複製されるように指定してもよい。予算ストレージクラスは、オブジェクト30の1つの複製が単一領域310に保存されるように指定してもよい。ローカルストレージクラスは、オブジェクト30が少なくとも3回、単一領域310に複製されるように指定してもよい。その他の実施形態において、ストレージサービスシステムは、その他の特徴を有するストレージクラスを定義する、又はクライアント50が、ストレージポリシー（例えば、ノードピッカー130に関して上記に記載されるような）の組み合わせを指定することにより、所定のオブジェクト30のストレージクラスをカスタマイズすることを可能にする。

#### 【0247】

上記に記載されるように、SLA特性は、複製の数及び複製が分散されるべき領域の数を越えて拡張してもよい。一実施形態において、特定のストレージクラスのSLA特性は、特定のストレージクラスに関連するオブジェクト30に対する期待処理待ち時間の指示を含む場合がある。例えば、1つのストレージクラスは、所定の費用で低期待処理待ち時間を指定してもよいが、別のストレージクラスは、より低コストで高期待処理待ち時間を指定してもよい。異なるレベルの期待処理待ち時間が、様々な方法で実装されてもよい。例えば、所定のコーディネータ120の観点から、いくつかのノード160は、ノード160と所定のコーディネータ120の近接ノード160で利用可能なリソースのレベル及び種類、ノード160の処理負荷、又はその他の関連した要因により、他より短いアクセス待ち時間を示すことができる。従って、所定のストレージクラスにより指定されるその他のSLA特性によりもたらされる制約に従い、一部の実施形態において、コーディネータ120及びノードピッカー130は、より短い期待処理待ち時間を指定するストレージクラスのオブジェクト30に対して、より短いアクセス待ち時間を示すノード160を選択するように設定されてもよい。その他の実施形態において、コーディネータ120は、オブジェクト30に関連するストレージクラスの期待処理待ち時間に従い、オブジェクト30に対するクライアントのアクセス要求の処理に優先順位をつけるように設定されてもよい。例えば、コーディネータ120は、低期待処理待ち時間を有するストレージクラスを支持し、処理にバイアスをかける一方、高期待処理待ち時間を有するストレージクラスの要求が、最終的には確実に完了するように設定されている、異なるキュー又はその他の処理制御、あるいはデータ構造を実装してもよい。

#### 【0248】

ストレージクラスは、オブジェクト30が最初にストレージサービスシステムに保存される時点で、クライアント50により指定されてもよい。あるいは、一部の実施形態において、クライアント50は、オブジェクト30がストレージサービスシステム内に存在する間のいずれの時点で、オブジェクト30に関するストレージクラスを変更してもよい。オブジェクト30が最初に保存される時点で、クライアント50によりストレージクラスが指定されていない場合、上記に記載されるような初期ストレージクラスが使用される。

上記に記載されるように、一部の実施形態において、オブジェクト 30 のストレージクラスは、オブジェクト 30 のキーに関連するキーマップ記録 148 内に保存されてもよい。そのような実施形態において、コーディネータ 120 及び / 又はレプリケータ 180 は、オブジェクト 30 の既存の複製を保存、複製、及び維持する際に、オブジェクト 30 のストレージクラスを考慮するように設定されてもよい。クライアント 50 は、異なるストレージクラスに関連するオブジェクト 30 に対して、異なる使用料を請求される場合があることが考慮される。例えば、高信頼性ストレージクラスは、一般的によりシステムリソースを使用する場合があります、一方、予算ストレージクラスは、より少ないリソースを使用してもよい。従って、所定のサイズのオブジェクト 30 において、クライアント 50 は、前者のストレージクラスを使用するオブジェクト 30 の保存に対してはより多く、後者に対してはより少なく請求される。

10

#### 【0249】

ストレージサービスシステム内のストレージクラスの動作方法の一実施形態を、図 26 に図示する。図示される実施形態において、動作は、クライアント 50 が特定のオブジェクト 30 と関連するストレージクラスを指定するブロック 2600 から始まる。続いて、ストレージクラスは、ストレージサービスシステム (ブロック 2602) 内において、特定のオブジェクト 30 と永続的に関連付けられる。例えば、ストレージクラスの指示は、クライアント 50 の代わりに、コーディネータ 120 によって、特定のオブジェクト 30 と関連するデータ構造内、例えばキーマップ記録 148 に保存される。オブジェクト 30 と関連するオブジェクトデータの状態は、そして指定されたストレージクラスの特性に従って設定される (ブロック 2604)。例えば、ストレージクラスが、領域 310 内のオブジェクトの複製の数及び / 又は分散に対して特定の要件を指定する場合、コーディネータ 120 及び / 又はレプリケータ 180 は、特定のオブジェクト 30 と関連するストレージシステムの結果的に生じる状態が、ストレージクラスの要件を満たすように、必要な複製を作成及び分散するように動作する。一部の実施形態において、レプリケータ 180 は、オブジェクト 30 に対するストレージクラス要件が長い間確実に維持されるように設定されてもよい。例えば、複製が機能しなくなった場合、レプリケータ 180 は障害を検出し、さらなる複製を作成するように設定されてもよい。

20

#### 【0250】

一部の実施形態において、所定のストレージクラスにより指定されるストレージ特性は、ビットストアノード 160 を介して利用可能な異なる種類のストレージリソースの識別を含む場合がある。例えば、一部の実施形態において、いくつかのビットストアノード 160 は、他より高性能なストレージデバイスを含む、又はそれぞれのビットストアノード 160 は、高性能及び低性能なデバイスの組み合わせを含む場合がある。そのような実施形態において、ストレージクラスは、そのクラスに関連するオブジェクト 30 に対して、どちらか一方のデバイスを使用するように指定してもよい。

30

#### 【0251】

##### 動的複製

上記に記載されるように、一部の実施形態において、ノードピッカー 130 は、特定のオブジェクト 30 の複製が書き込まれるべき特定のビットストアノード 160 を識別する書き込みプランを作成するように設定されてもよい。そのような書き込みプランは、一度書き込みプランが、例えばコーディネータ 120 などによって導入されると、例えば特定のオブジェクト 30 に関する様々な書き込みポリシーが満たされるように作成されてもよい。例えば、書き込みプランにより指定されるノード 160 の数は、特定のオブジェクト 30 の複製の最小要求数、複製を分散する異なる領域 310 の最小数、又はあらゆるその他のストレージポリシー配慮に従い、決定されてもよい。

40

#### 【0252】

一部の実施形態において、ノードピッカー 130 は、ノード 160 の現行状態を考慮することなく、予測可能な手順に従い、ノード 160 を連続して選択する静的方法にて、書き込みプランを作成するように設定されてもよい。例えば、ノードピッカー 130 は、文

50

字列複製のためにノード１６０の同一組を連続して選択する、又はラウンドロビン法で汎山のノード１６０を回転してもよい。しかし、大規模な実装において、ストレージサービスシステムは、様々な時点において大幅に異なる状態で動作し得る、多くのノード１６０を含む場合がある。例えば、いくつかのノード１６０は動作不能であり、その他は動作しているが要求アクティビティにより飽和状態であるか、又は使用可能なリソースが少ない、そしてその他は比較的遊休状態又は十分な使用可能なリソースを有してもよい。

#### 【０２５３】

さらに、異なるノード１６０は、それぞれの任意のコーディネータ１２０又はノードピッカー１３０から、異なるレベルの通信費用を提示してもよい。例えば、コーディネータ１２０と同一の領域３１０又はデータセンタ３００内に存在するノード１６０は、ローカルで応答時間の短いネットワーク接続を介して、アクセス可能であってもよい。一方、コーディネータ１２０とは異なる領域３１０又はデータセンタ３００に存在するノード１６０は、ローカルノード１６０より大幅に長い待ち時間を示すことができる。さらに、一部の実施形態において、領域３１０間又はデータセンタ３００間の通信は、ローカル通信とは異なる経済費用モデルの通信ネットワークにおいて行われる場合がある。例えば、領域３１０内での通信は、データ転送において、使用量に基づく請求のない、十分なバンド幅を有するプライベートローカルエリアネットワーク（ＬＡＮ）において行われる場合がある。一方、データセンタ３００間の通信は、専用通信設備、公共インターネット、プライベート広域ネットワーク（ＷＡＮ）設備、又はその他の長距離通信ネットワークなどの設備において行われる場合がある。これらの設備は、一般的にＬＡＮ設備よりバンド幅に制限がある場合があり、一部の例において、ＬＡＮ通信には適用されない利用費用（例えば、ピーク又は合計バンド幅使用量に基づく）が第三者により請求されてもよい。

#### 【０２５４】

様々なノード１６０の両方の動作状態及びこれらのノードの通信費用は、長期に渡り、変動してもよい。例えば、ある時点において動作している、又は遊休状態のノード１６０は、後に動作不能又は使用中になる場合があり、またその逆も同様である。同様に、待ち時間及び／又は経済費用などの通信費用は、ある期間において高くなり、その他では低くなる（例えば、ピーク時対ピーク時以外の利用量）場合がある。この変動のため、ある時点において効率的で低コストな書き込みプランは、大幅に効率が低く、高コストである、また別の時点（例えば、書き込みプランに指定されているノード１６０が使用中になる、通信が遅くなる、又は動作不能になる場合）において、実行不可能でさえある場合もある。

#### 【０２５５】

従って、一部の実施形態において、ノードピッカー１３０は、任意のオブジェクト３０の書き込み複製に対する任意の書き込みプランを、ノード１６０の現行状態情報に従い、動的に決定するように設定されてもよい。一般的に言うと、動的に決定される書き込みプランは、ノード１６０の監視可能な動的状態情報を考慮してもよい。すなわち、動的に決定される書き込みプランは、時間とともに変化するノードの状態情報の機能として作成されてもよい。従って、任意のオブジェクト３０の動的に決定される書き込みプランは、ノード１６０の状態から独立して決定されてもよい、静的に作成される書き込みプランと比較し、ノード１６０の基礎状態情報により、時間とともにそれ自身を変化してもよい。

#### 【０２５６】

上記に記載されるように、多くの異なる種類の状態情報は、書き込みプランの動的作成に考慮されてもよい。一般的に、ノード１６０の状態情報は、任意のノード１６０に関する状態情報、及びそれを介して任意のノード１６０にアクセス可能な通信リソースに関する状態情報（例えば、ネットワークリソース）を含む場合がある。様々な実施形態において、任意のノード１６０に関する状態情報は、任意のノード１６０（又はノードに関連するアプリケーションインスタンス）がＯＫであるかどうか、隔離されているか、又は上記に記載されるようなＤＦＤＤ１１０により示されるその他の動作状態などの、任意のノードの動作状態を含む場合がある。任意のノード１６０に関する状態情報は、動作状態情報

より詳細な任意のノード 160 の行動を示す場合がある負荷状態情報も含む場合がある。例えば、様々な実施形態において、負荷状態情報は、任意のノード 160 に対応するプロセッサ利用レベル、メモリ利用レベル、ストレージデバイス利用レベル、ストレージデバイス入力/出力バンド幅利用レベル、又はネットワークインターフェースのバンド幅利用レベル、又はノード動作のいずれのその他の計測可能な局面を示すことができる。上記に記載されるように、一部の実施形態において、動作状態情報に加え、負荷状態情報は、DFDD110 を介して入手可能であってもよい。

#### 【0257】

ネットワーク費用情報とも称される場合がある通信リソース状態情報は、任意のノード 160 への 1 つ以上の通信経路の状態に関するいずれの適切な情報を含む場合がある。様々な実施形態において、ネットワーク費用情報は、任意のノード 160 への及び/又はからのメッセージの伝達に関連するネットワーク通信待ち時間を示す場合があり、時間（秒、ミリ秒など）、ネットワークホップ数（例えば、メッセージを伝達するためのルーティングステップ数）、又は別の適切な測定基準により表されてもよい。一実施形態において、ネットワーク費用情報は、任意のノード 160 との通信に利用可能な利用可能バンド幅（例えば、データ転送率）の指示を含む場合がある。別の実施形態において、ネットワーク費用情報は、任意のノード 160 とのネットワーク通信に係る経済費用の指示を含む場合がある。例えば、そのような費用は、ある量のデータの転送又は受信に対して請求される率、又はいずれのその他のネットワーク通信に適切な費用又は率モデルにより表されてもよい。

#### 【0258】

ノードピッカー 130 は、通常、オブジェクト 30 の書き込みプランの動的な決定において、ノード 160 の状態情報のいずれの適切な機能を使用してもよい。一部の実施形態において、ノードピッカー 130 により導入されるストレージポリシー（前述のストレージポリシーの追加又は代替であってもよい）は、特定のオブジェクト 30 に対する書き込みプランに含めることが可能である、ノード 160 を制約する状態情報に対するガイドライン又は要件を指定してもよい。様々な実施形態において、これらのポリシーは、オブジェクト 30 の特定の組（例えば、特定のストレージクラス又はパケットに含まれており、共通鍵プレフィックスを有する、又は組のメンバーである则表示されている、オブジェクト）に対して全体的に（例えば、すべてのオブジェクト 30 に対して）、又は個別のオブジェクト 30（例えば、オブジェクト 30 に関する特定のポリシーを指定するクライアントへの応答として）に対して、あるいはいずれの適切なこれらの組み合わせで適用されてもよい。一実施例において、特定のストレージクラスは、いくつかの最小数の複製は、いくつかの最大通信待ち時間以下を示すことを要求するストレージポリシーを指定してもよい。同様に、このストレージクラスのオブジェクト 30 への書き込みプランの作成において、ノードピッカー 130 は、指定される最大通信待ち時間を満たすかどうかに従い、少なくともいくつかのノード 160 を選択するように設定されてもよい。

#### 【0259】

一部の実施形態において、ノードピッカー 130 は、ノード状態情報における様々な種類の最適化に従い、書き込みプランを作成するようにも設定されてもよい。例えば、特定の最大ネットワーク費用又はその他の書き込みプランに関連する費用を指定する代わりに、ストレージポリシーは、特定の時間において利用可能なリソース間で費用が最小限化されることを指定してもよい。同様に、ノードピッカー 130 は、例えば、低ネットワーク通信又はその他の関連費用を有するノード 160 を選択することによって、書き込みプランに関連する 1 つ以上の費用を最小化するように設定されてもよい。一部の実施形態において、そのような最小化は、その他のノード状態情報要件を指定するその他のストレージポリシーなどの、その他の制約の存在下において起こる場合がある。

#### 【0260】

さらに、一部の実施形態において、いくつかのノード状態情報は、予測可能な形で時間とともに変化することに留意されたい。例えば、データセンタ 300 間のネットワーク通

10

20

30

40

50

信に関連するバンド幅費用は、明確な料金表に従い、変化する。一部の実施形態において、書き込みプランに関連する費用を最小化することは、特定の時間帯に関連する費用に基づき、書き込みプランのすべて又は一部が実行されるべき時間帯を識別することを含む場合がある。例えば、ノードピッカー 130 は、遠隔データセンタ 300 と通信するためのバンド幅が、将来のある時点において現時点より安価になることを判断する場合があり、さらに遠隔データセンタ 300 に存在するノード 160 を含む書き込みプランの費用は、遠隔データセンタを対象とする少なくともそれらのストレージ作業を、特定の将来時に実行することで最小化できる可能性があることを判断する。このプロセスの 1 つの起こりうる結果は、ノードピッカー 130 により作成される書き込みプランが、任意のオブジェクト 30 のいくつか（又は場合によりすべて）の複製の作成は、特定の将来時まで延期されるべきであるということを示す場合があることである。

10

#### 【0261】

多くの異なるストレージポリシーを特定のオブジェクト 30 に適用することは可能である。さらに、ある場合において、特定のオブジェクト 30 に関連する各ストレージポリシーを満たす 1 つの書き込みプランを作成することは不可能である場合がある。例えば、特定のオブジェクト 30 に関連するストレージポリシーは、複製の最小数を、最小数の異なる領域 310 に渡り、保存し、配布することを指定してもよい。しかし、オブジェクト 30 に対して、書き込みプランが作成された時点では、ノードピッカー 130 が実行中の領域 310 は、一時的通信障害により、その他の領域 310 から一時的に隔離されてもよい。従って、対応するストレージポリシーを満たしつつ、複製をその他の領域 310 に正常に配布することは、少なくとも一時的に不可能である場合がある。

20

#### 【0262】

一実施形態において、ノードピッカー 130 は、オブジェクト 30 に対する書き込みプランを、書き込みプランが満たすストレージポリシーの数を最大化するという原則において、動的に決定するように設定されてもよい。次善条件の存在下において、これは、ストレージポリシーを満たすために「最善努力」を示す書き込みプランをもたらす場合がある。例えば、今記載した特定のシナリオにおいて、通信障害のため、領域別ポリシーは満たせない場合があるが、最小複製ポリシーは、特定のオブジェクト 310 の要求される最小数の複製を、ローカル領域 310 内に保存することで、満たされてもよい。一部の実施形態において、ストレージポリシーの最大化は、様々な抑制化で行われる場合がある。例えば、それらが満たされない場合、書き込みプランは決定されず、クライアントのオブジェクトを保存するという要求は、失敗する可能性があるというような、いくつかのストレージポリシーは、義務として認識されてもよい。その他のストレージポリシーは、選好関係又は重み付けを有する場合があり、例えば、最大化プロセスにおいて、低い選好ストレージポリシーの中から、より高い選好ストレージポリシーが選択されてもよい。別の実施形態において、ストレージポリシーの選択は、結果として生じるストレージプランにより満たされるストレージポリシーの数の代わりに、又はそれへの追加として、結果として生じるストレージプランの合計重み（満たされるストレージポリシーの重みに基づき決定される）を最大化することにより実行されてもよい。

30

#### 【0263】

オブジェクト 30 に対する書き込みプランを動的に決定するための様々な技術は、オブジェクト 30 が最初に保存される際、単独で実行する必要はないことに留意されたい。上記に記載されるように、一部の実施形態において、レプリケータ 180 は、オブジェクト 30 の複製がアクセス可能であるかどうかを決定するために、オブジェクト 30 に対応するキーマップエントリ 144 を調査するように設定されてもよい。特定のオブジェクト 30 のいずれの複製もアクセス可能でない場合、レプリケータ 180 は、さらなる複製を作成するために使用されてもよい、新しい書き込みプランをノードピッカー 130 から要求するように設定されてもよい。新しい書き込みプランは、上記に記載される技術のいずれの適切な組み合わせを使用するノードピッカー 130 により、動的に決定されてもよい。さらに、一部の実施形態において、レプリケータ 180 は、オブジェクト 30 の様々なスト

40

50

レージポリシーへの順守をより広く監視するように設定されてもよい。例えば、レプリケータ 180 は、オブジェクト 30 の既存の複製の組が、最小複製ポリシー、又はいずれのその他の適切なポリシーの組に加え、領域別ポリシーを満たすかどうかを判断するように設定されてもよい。そのような一実施形態において、レプリケータ 180 が、特定のオブジェクト 30 の既存の複製により満たされているポリシーの数が、基準値より少ないと判断する場合、レプリケータ 180 は、上記に記載されるように、満たされるストレージポリシーが最大になるように動的に決定される新しいストレージプランをノードピッカー 130 から要求してもよい。別の実施形態において、レプリケータ 180 は、特定の義務ストレージポリシーが満たされないという判断に伴い、又は満たされるストレージポリシーの合計重みが基準値を下回るという判断に伴い、新しいストレージプランを要求してもよい。

10

#### 【0264】

図 27 は、ビットストアノード 160 の現行状態情報に従い、データオブジェクトの 1 つ以上の複製を保存するための書き込みプランを動的に決定する方法の一実施形態を図示する。図示される実施形態において、操作は、任意のオブジェクト 30 を保存するというクライアントの要求が受け取られるブロック 2700 から始まる。一実施形態において、そのような要求は、上記に詳細が記載されるウェブサービスインターフェース 100 を介して、ウェブサービスプロトコルに従い受け取られる場合がある。

#### 【0265】

続いて、ビットストアノード 160 の現行状態情報に従い、任意のオブジェクト 30 の複製を保存するための書き込みプランが動的に決定される（ブロック 2702）。例えば、ノードピッカー 130 は、任意のオブジェクト 30 に適用されてもよい様々なストレージポリシーに従い、書き込みプランを決定するように設定されている場合があり、当該ポリシーは、ノードの動作状態、ノードの負荷状態情報、ネットワーク通信費用、又は上記に詳細に記載されるような、いずれのその他の適切な状態情報などの、いずれの適切な現行状態情報を考慮する。さらに、上記に記載されるように、一部の実施形態において、動的な書き込みプランの決定は、書き込みプランに関連する費用を最小化すること、又は書き込みプランにより満たされるストレージポリシーの重み又は数を最大化することなどによる、状態情報又はストレージポリシーに関する最適化を含む場合がある。

20

#### 【0266】

任意のオブジェクト 30 の複製は、そして動的に決定された書き込みプランに従い、1 つ以上のビットストアノード 160 に保存される（ブロック 2704）。例えば、コーディネータ 120 は、上記に記載されるように、書き込みプランに指定されるそれぞれのビットストアノード 160 に対する操作を設置するビットストアオブジェクトを作成するように設定されてもよい。一部の実施形態において、書き込みプランのいくつかのストレージ操作は、上記に記載されるように、その他の操作とは異なる時間に実行されてもよい。

30

#### 【0267】

前述のように、一部の実施形態において、書き込みプランは、1 つ以上の複製がビットストアノード 160 に渡り既に保存されているオブジェクト 30 に対して動的に決定されてもよい。図 28 は、そのような方法の一実施形態を図示する。図示される実施形態において、操作は、任意のオブジェクト 30 の 1 つ以上の既存の複製が調査されるブロック 2800 から始まる。例えば、上記に記載されるように、レプリケータ 180 の一実施形態は、任意のオブジェクト 30 の既存の複製がアクセス可能であるか、及び / 又は任意のオブジェクト 30 の既存の複製がどの程度オブジェクトに関連するストレージポリシーを満たすかどうかを判断するように設定されてもよい。

40

#### 【0268】

任意のオブジェクト 30 の複製の調査を受けて、1 つ以上のさらなる複製を作成する必要があると判断されてもよい（ブロック 2802）。例えば、既存の複製は、故障しているか、又は別の理由でアクセス不可となり、複製の最小数を下回っている場合がある。あるいは、既存の複製の状態は、1 つ以上のストレージポリシーに対して不十分である場合

50

がある。続いて、任意のオブジェクト 30 のさらなる複製を保存するための書き込みプランは、ビットストアノード 160 の現行状態情報に従い、動的に決定される（ブロック 2804）。そのような書き込みプランは、前述と同様に、又はいずれの適切なその変更に従い、決定されてもよい。一部の実施形態において、任意のオブジェクト 30 のさらなる複製を作成する必要がないと判断される場合、書き込みプランは決定されない場合があることに留意されたい。

#### 【0269】

任意のオブジェクト 30 の複製は、そして動的に決定された書き込みプラン（ブロック 2806）に従い、1つ以上のビットストアノード 160 に保存される。例えば、レプリケータ 180 は、上記に記載されるように、書き込みプランに指定されるそれぞれのビットストアノード 160 に対して操作を設置するビットストアオブジェクトを作成するように設定されているか、又は単に任意のオブジェクト 30 の既存の複製を保存する1つ以上のノード 160 を、書き込みプランに指定されるノード 160 の1つ以上に直接それらの複製をコピーしてもよい。

#### 【0270】

##### 例示的なコンピュータシステムの実施形態

一部の実施形態において、上記に記載されるいずれの方法又は技術は、プログラムのインストラクション及びコンピュータアクセス可能な媒体を介して保存、又は伝達することが可能なデータとして実装されてもよい。そのような方法又は技術は、例えば、ストレージクライアント 50、ウェブサービスプラットフォーム 100、DFDD 110、コーディネータ 120、ノードピッカー 130、キーマップインスタンス 140、ビットストアノード 160、レプリケータ 180、及び/又はレプリケータキーマップ 190 を含む場合があるが、制限はない。そのような方法又は技術は、図 6～9、13～15、20～22、及び 25～28 に図示されているいずれの方法ならびにその適切な変更をさらに含む場合がある。そのようなプログラムのインストラクションは、上記に記載される特定の方法又は一部の方法などの特定のコンピュータ機能を実行し、より一般的なオペレーティングシステムの機能性、アプリケーションの機能性、及び/又はいずれのその他の適切な機能を提供するために、実行されてもよい。一部の実施形態において、上記に記載されるコンポーネント又は方法は、異なるその他の実施形態において、これらの示されているものより少ないエンティティに組み込まれる、又は機能性はコンポーネント又は方法に渡り、上記に記載される分割とは異なる形で分割されてもよい。

#### 【0271】

コンピュータアクセス可能な媒体を含むコンピュータシステムの例示的な一実施形態を図 29 に図示する。そのようなシステムは、ノードとも称される場合がある。前述のように、一実施形態において、上記に記載されるいずれの様々なストレージシステムコンポーネントは、多くのノードを渡り分散されている場合があり、そのような任意のコンポーネントは、1つ以上のノードにより実装されている、又はいくつかのノードに渡り分割されてもよい。一部の実施形態において、ノードは単一ストレージサービスシステムコンポーネントの機能を排他的に実装してもよいが、その他の実施形態において、ノードは、いくつかの異なるシステムコンポーネントのすべて、又は部分的な機能性を実装してもよい。図示される実施形態において、コンピュータシステム 2900 は、入力/出力（I/O）インターフェース 2930 を介してシステムメモリ 2920 と対になっている1つ以上のプロセッサ 2910 を含む。コンピュータシステム 2900 は、I/O インターフェース 2930 と対になっているネットワークインターフェース 2940 をさらに含む。

#### 【0272】

様々な実施形態において、コンピュータシステム 2900 は、1つのプロセッサ 2910 を含むユニプロセッサシステム、又はいくつかのプロセッサ 2910（例えば、2つ、4つ、8つ、又は別の適切な数）を含むマルチプロセッサシステムであってもよい。プロセッサ 2910 は、命令を実行する能力を有する、任意の適切なプロセッサであってもよい。例えば、様々な実施形態において、プロセッサ 2910 は、x86、PowerPC

、S P A R C、又はM I P S I S A、あるいはその他の適切なI S Aなどの様々な命令組アーキテクチャ(I S A)を実装する汎用、又は内蔵プロセッサであってもよい。マルチプロセッサシステムにおいて、それぞれのプロセッサ2910は、必ずではないが、通常、同一I S Aを実装する。

#### 【0273】

システムメモリ2920は、プロセッサ2910によりアクセス可能に命令及びデータを保存するように設定されてもよい。様々な実施形態において、システムメモリ2920は、静的ランダムアクセスメモリ(S R A M)、同期動的R A M(S D R A M)、非揮発性/フラッシュ型メモリ、又はいずれのその他の種類のメモリなど、いずれの適切なメモリ技術を使用し、実装されてもよい。図示される実施形態において、上記に詳細を記載する、いずれのこれらのストレージサービスシステムコンポーネント及びその他の機能などの所望の機能を実装するプログラムのインストラクション及びデータは、システムメモリ2920内にコード2925として保存されることを示す。

#### 【0274】

一実施形態においてI/Oインターフェース2930は、ネットワークインターフェース2940又はその他の周辺インターフェースを含む、プロセッサ2910、システムメモリ2920、及びデバイス内のいずれの周辺装置間のI/Oトラフィックを調整するように設定されてもよい。一部の実施形態において、I/Oインターフェース2930は、いずれの必要なプロトコル、タイミング調節、又はその他のあるコンポーネント(例えば、システムメモリ2920)からのデータ信号を別のコンポーネント(例えば、プロセッサ2910)での使用に適切なフォーマットに変換するためのデータ変換を実行してもよい。一部の実施形態において、I/Oインターフェース2930は、例えば、周辺装置相互接続(P C I)バス標準又はユニバーサルシリアルバス(U S B)標準の異型などの様々な種類の周辺装置のバスを通じて接続されたデバイスのサポートを含む場合がある。一部の実施形態において、I/Oインターフェース2930の機能は、例えば、ノースブリッジ及びサウスブリッジなど、2つ以上の分離したコンポーネントに分割されてもよい。また、一部の実施形態において、システムメモリ2920に対するインターフェースなどの、I/Oインターフェース2930の機能性の一部又はすべては、プロセッサ2910に直接組み込まれている場合がある。

#### 【0275】

ネットワークインターフェース2940は、コンピュータシステム2900とネットワークに接続されたその他のデバイス、例えばその他のコンピュータシステムなどとのデータ交換が可能になるように設定されてもよい。様々な実施形態において、ネットワークインターフェース2940は、有線又は無線の一般的なデータネットワーク、例えば、イーサネット(登録商標)ワークのいずれの適切な種類、アナログ音声ネットワーク又はデジタルファイバー通信ネットワークなどの電気通信/電話ネットワーク、ファイバーチャネルS A Nなどのストレージ領域ネットワーク、又はその他のいずれの適切な種類のネットワーク又はプロトコルなどを介して、通信をサポートしてもよい。

#### 【0276】

一部の実施形態において、システムメモリ2920は、上記に記載されるように、プログラムのインストラクション及びデータを保存するように設定されたコンピュータアクセス可能媒体の一実施形態であってもよい。しかし、その他の実施形態において、プログラムのインストラクション及び/又はデータは、異なる種類のコンピュータアクセス可能媒体から受け取る、へ送る、上に保存してもよい。一般的に、コンピュータアクセス可能媒体は、例えば、I/Oインターフェース2930を介してコンピュータシステム2900と対になっているディスク又はC D / D V D - R O Mなどの磁気又は光学媒体などのストレージ媒体あるいはメモリ媒体を含む場合がある。またコンピュータアクセス可能媒体は、コンピュータシステム2900の一部の実施形態に、システムメモリ2920又は別の種類のメモリとして含まれている場合がある、R A M(例えばS D R A M、D D R S D R A M、R D R A M、S R A Mなど)、R O M、などのいずれの揮発性又は非揮発性媒体



を含む場合もある。コンピュータアクセス可能媒体を介して保存されたプログラムのインストール及びデータは、転送媒体、又は電気、電磁気、又はネットワーク及び/又はネットワークインターフェース2940を介して実装されている無線リンクなどの通信媒体を介して伝送される場合があるデジタル信号などの信号により、転送されてもよい。

【0277】

上記の実施形態は、相当な詳細を記載したが、多数の変更及び修正は、一度上記の開示が完全に理解されると、当技術分野に精通する者により明らかになるであろう。以下の特許請求の範囲は、すべてのそのような変更及び修正を包含することが意図される。

【図面の簡単な説明】

【0278】

10

【図1】ウェブサービスとしてユーザにストレージを提示するためのストレージモデルの一実施形態を示すブロック図である。

【図2】ストレージサービスシステムのアーキテクチャの一実施形態を示すブロック図である。

【図3】ストレージシステムコンポーネントの物理的配置の一実施形態を示すブロック図である。

【図4】ストレージノードの一実施形態を示すブロック図である。

【図5】ストレージノード内のデータオブジェクトを体系化するように構成されるデータストラクチャの一実施形態を示すブロック図である。

【図6】オブジェクト取得操作を実行する方法の一実施形態を示す工程図である。

20

【図7】オブジェクト格納操作を実行する方法の一実施形態を示す工程図である。

【図8】オブジェクト解放操作を実行する方法の一実施形態を示す工程図である。

【図9】オブジェクトストレージ空間を再圧縮する方法の一実施形態を示す工程図である。

【図10】キーマップインスタンスデータ構造の組の一実施形態を示すブロック図である。

【図11A】キーマップインスタンスの階層的実装の一実施形態を示す。

【図11B】キーマップインスタンスの階層的実装の一実施形態を示す。

【図11C】キーマップインスタンスの階層的実装の一実施形態を示す。

【図11D】キーマップインスタンスの階層的実装の一実施形態を示す。

30

【図12】キーマップインスタンス内の階層的層間における関係を要約するブロック図である。

【図13】キーマップエントリ格納操作を実行する方法の一実施形態を示す工程図である。

【図14】キーマップエントリが取得操作を実行する方法の一実施形態を示す工程図である。

【図15A】更新伝播を使用してキーマップインスタンスを同期化する方法の一実施形態を示す工程図である。

【図15B】アンチエントロピープロトコルを使用してキーマップインスタンスを同期化する方法の一実施形態を示す工程図である。

40

【図16】レプリケータキーマップエントリの一実施形態を示すブロック図である。

【図17】不均衡インデックスデータ構造の一実施形態を示す。

【図18】不均衡データ構造において使用するためのインデックスノードの一実施形態を示す。

【図19】層別インデックスデータ構造の一実施形態を示す。

【図20】不均衡インデックスデータ構造をトラバースするための方法の一実施形態を示す工程図である。

【図21】FINGERPRINTアンチエントロピープロトコルメッセージを処理する方法の一実施形態を示す工程図である。

【図22】FILTERアンチエントロピープロトコルメッセージを処理するための方法

50

の一実施形態を示す工程図である。

【図 2 3】発見及び障害検出デモン (DFDD) の一実施形態を示す。

【図 2 4】DFDD インスタンスによって維持してもよい大域操作状態マシンの一実施形態を示す。

【図 2 5】ゴシッププロトコルに従ったDFDD インスタンスを同期化するための方法の一実施形態を示す工程図である。

【図 2 6】ストレージサービスシステム内のストレージクラスの操作方法の一実施形態を示す工程図である。

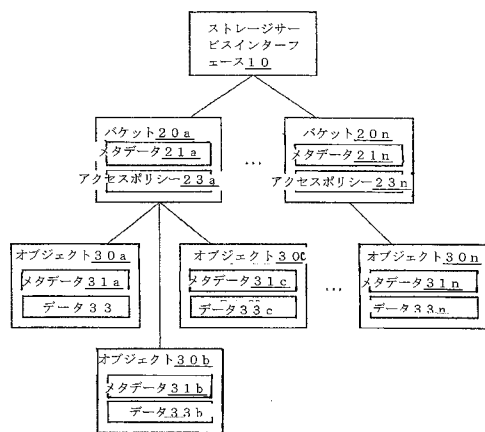
【図 2 7】ストレージノードの現在の情報に従い、データオブジェクトの 1 つ以上の複製を保存するための書き込みプランを動的に決定する方法の一実施形態を示す工程図である

10

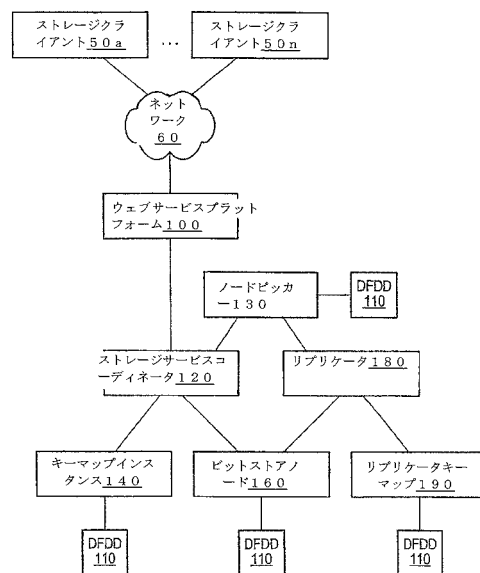
。【図 2 8】1 つ以上の複製が既にストレージノード間に保存されているオブジェクトに係る書き込みプランを動的に決定する一実施形態を示す工程図である。

【図 2 9】コンピュータシステムの例示的な実施形態を示す工程図である。

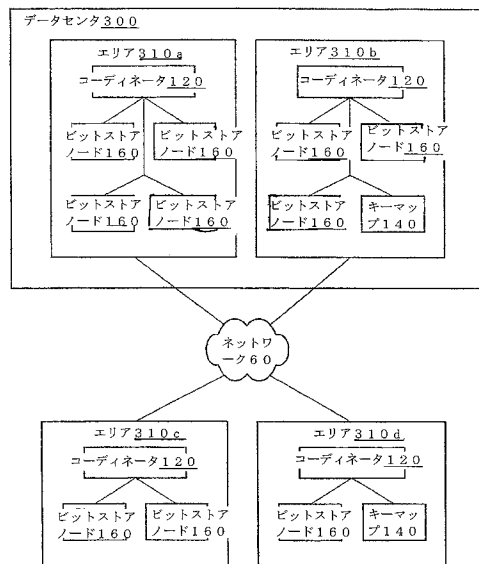
【図 1】



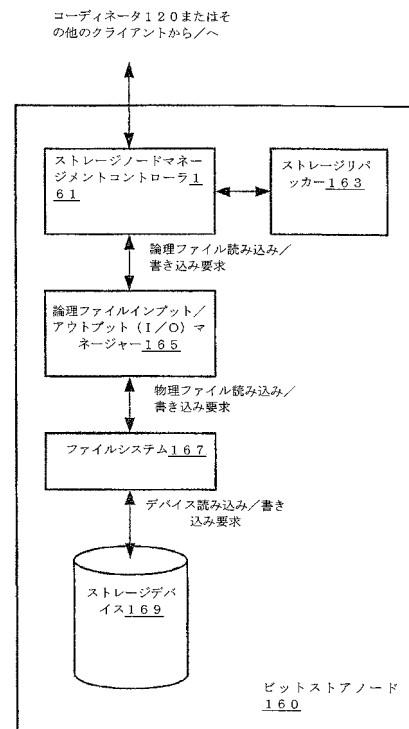
【図 2】



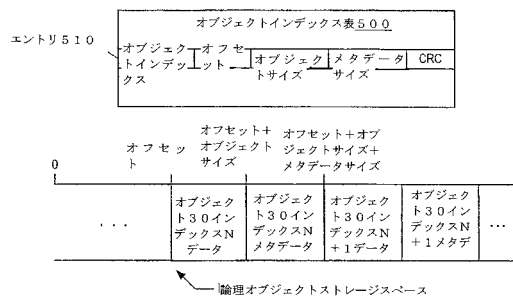
【図 3】



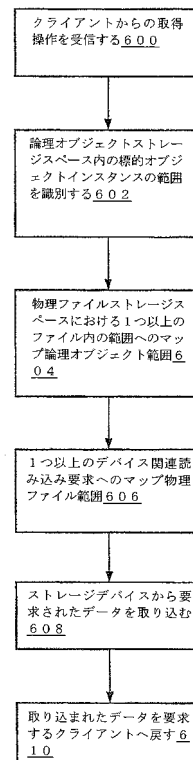
【図 4】



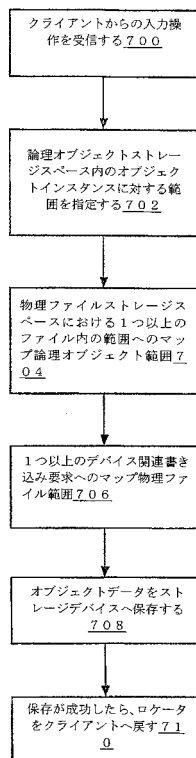
【図 5】



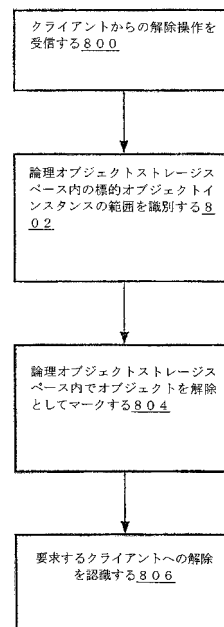
【図 6】



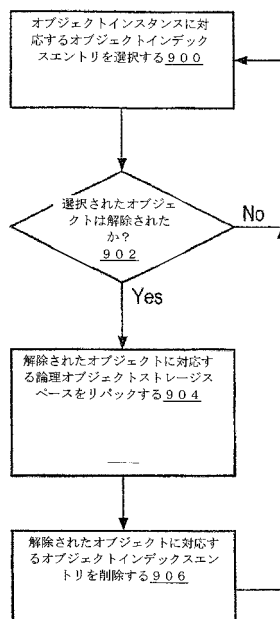
【図 7】



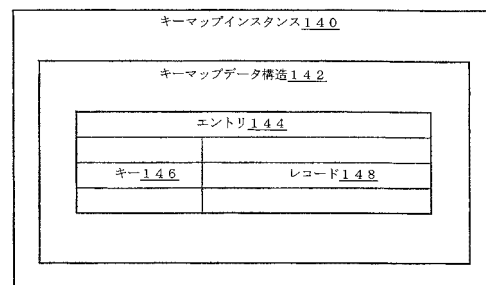
【図 8】



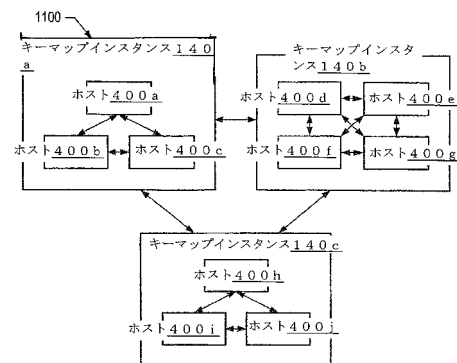
【図 9】



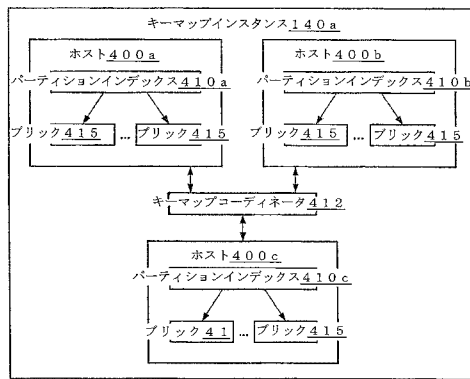
【図 10】



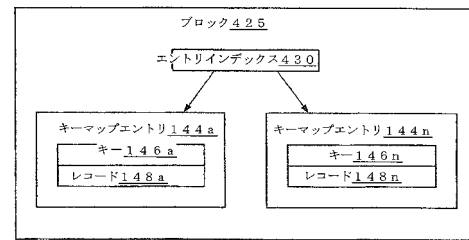
【図 11 A】



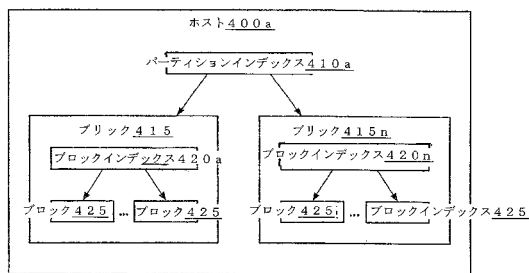
【図 11 B】



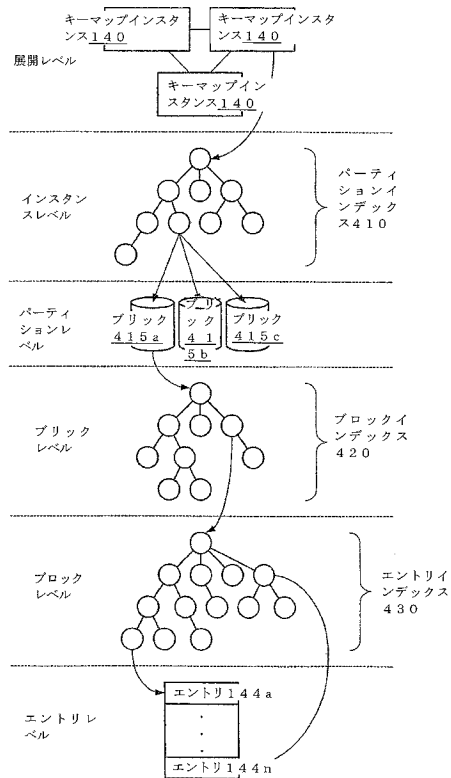
【図 11 D】



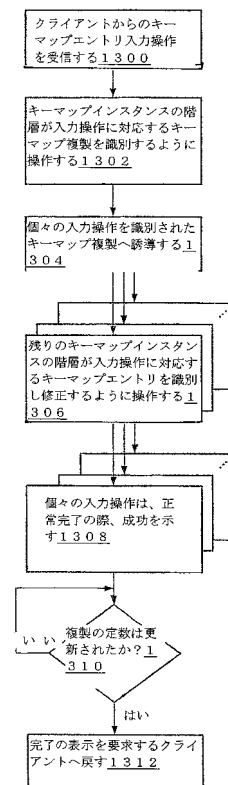
【図 11 C】



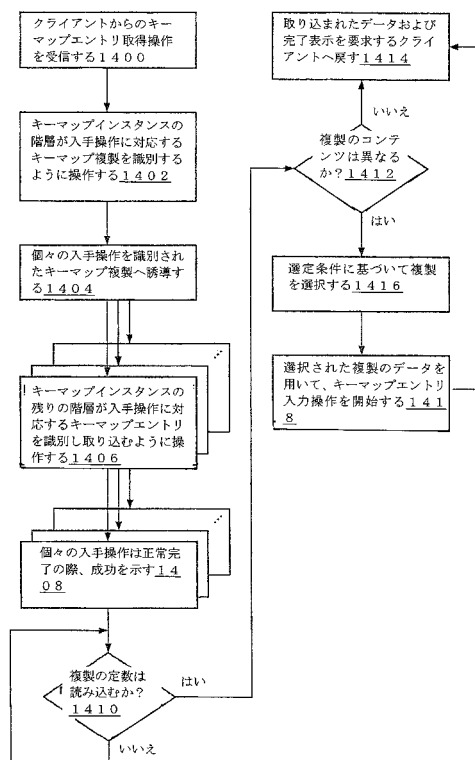
【図 12】



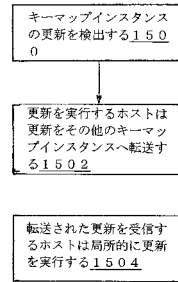
【図 13】



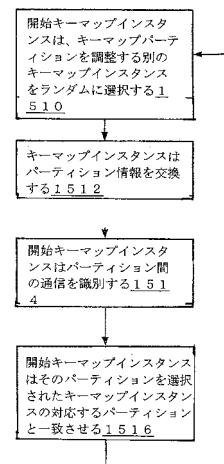
【図14】



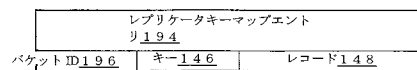
【図15A】



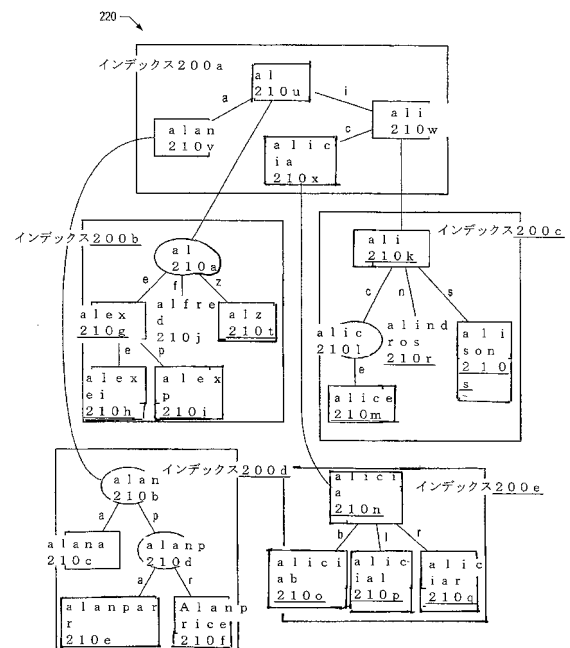
【図15B】



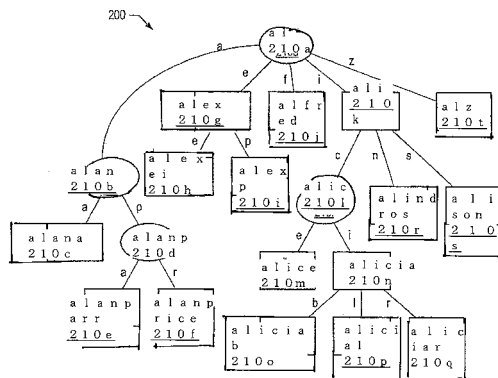
【図16】



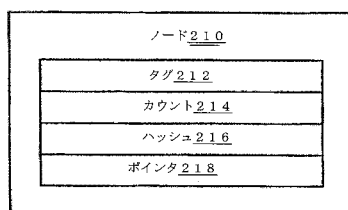
【図19】



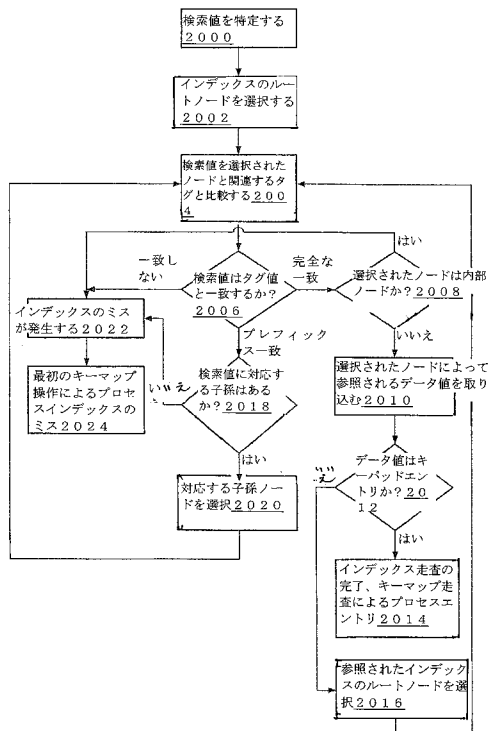
【図17】



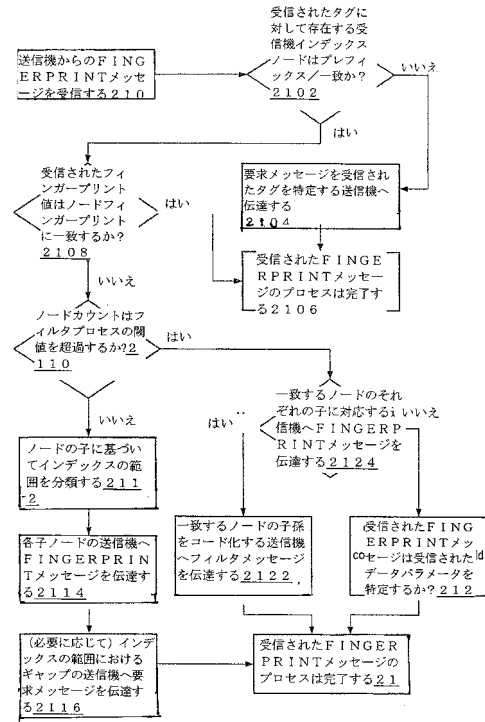
【図18】



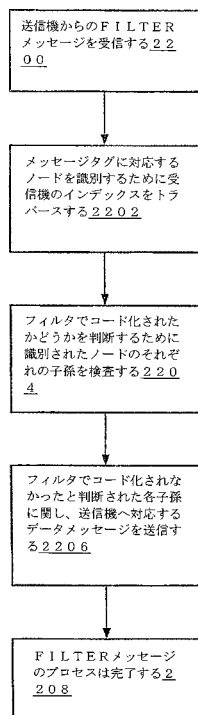
【図 20】



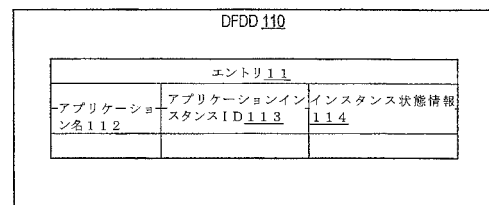
【図 21】



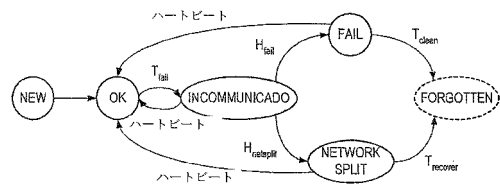
【図 22】



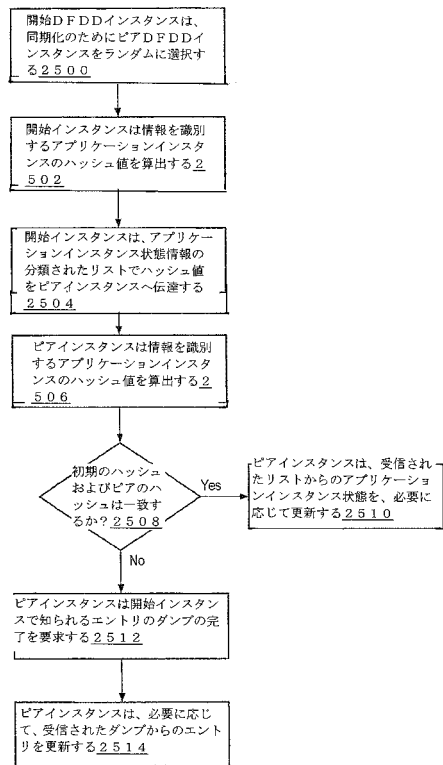
【図 23】



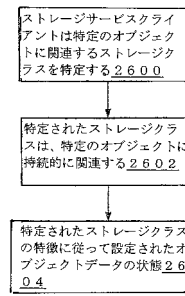
【図 24】



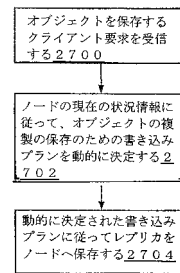
【図 25】



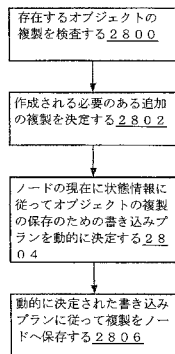
【図 26】



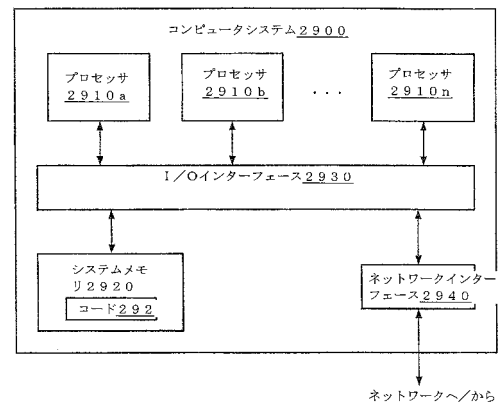
【図 27】



【図 28】



【図 29】





## フロントページの続き

- (72)発明者 アトラス, アラン・ビイ  
アメリカ合衆国・98144・ワシントン州・シアトル・12ティエイチ アベニュー サウス・1  
200・スイート 1200
- (72)発明者 バース, デイビッド・エム  
アメリカ合衆国・98144・ワシントン州・シアトル・12ティエイチ アベニュー サウス・1  
200・スイート 1200
- (72)発明者 コーミー, ジョン・デイビッド  
アメリカ合衆国・98144・ワシントン州・シアトル・12ティエイチ アベニュー サウス・1  
200・スイート 1200
- (72)発明者 フィッシュマン, エイミ・ケイ  
アメリカ合衆国・98144・ワシントン州・シアトル・12ティエイチ アベニュー サウス・1  
200・スイート 1200
- (72)発明者 ソレンソン, ジェイムズ・クリストファ・サード  
アメリカ合衆国・98144・ワシントン州・シアトル・12ティエイチ アベニュー サウス・1  
200・スイート 1200
- (72)発明者 ワグナー, エリック・エム  
アメリカ合衆国・98144・ワシントン州・シアトル・12ティエイチ アベニュー サウス・1  
200・スイート 1200

審査官 桜井 茂行

- (56)参考文献 国際公開第00/026782(WO, A1)  
特開2002-055869(JP, A)  
欧州特許出願公開第01160692(EP, A1)

(58)調査した分野(Int.Cl., DB名)

G06F 12/00  
G06F 17/30  
G06F 3/06